

# Object-based verification of a prototype Warn-on-Forecast system

Patrick S. Skinner<sup>\*1,2</sup>, Dustan M. Wheatley, Kent H. Knopfmeier<sup>1,2</sup>, Anthony E. Reinhart<sup>1,2</sup>,  
Jessica J. Choate<sup>1,2</sup>, Thomas A. Jones<sup>1,2</sup>, Gerald J. Creager<sup>1,2</sup>, David C. Dowell<sup>3</sup>, Curtis R.  
Alexander<sup>3</sup>, Therese T. Ladwig<sup>3,4</sup>, Louis J. Wicker<sup>2</sup>, Pamela L. Heinselman<sup>2</sup>, Patrick Minnis<sup>5</sup>,  
Rabindra Palikonda<sup>6</sup>

<sup>1</sup>*Cooperative Institute for Mesoscale Meteorological Studies, University of Oklahoma, Norman,  
Oklahoma*

<sup>2</sup>*NOAA/OAR/National Severe Storms Laboratory, Norman, Oklahoma*

<sup>3</sup>*NOAA/OAR/Earth System Research Laboratory, Boulder, Colorado*

<sup>4</sup>*Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder,  
Colorado*

<sup>5</sup>*NASA Langley Research Center, Hampton, Virginia*

<sup>6</sup>*Science Systems and Applications Inc., Hampton, Virginia*

<sup>\*</sup>*Corresponding author address:* Cooperative Institute for Mesoscale Meteorological Studies, University of Oklahoma, 120 David L. Boren Blvd. Norman, OK 73072.

E-mail: patrick.skinner@noaa.gov

## ABSTRACT

17 An object-based verification methodology for the NSSL Experimental  
18 Warn-on-Forecast System for ensembles (NEWS-e) has been developed and  
19 applied to 32 cases between December 2015 and June 2017. NEWS-e fore-  
20 cast objects of composite reflectivity and 30-minute rotation tracks of updraft  
21 helicity are matched to corresponding objects in Multi-Radar Multi-Sensor  
22 data on space and time scales typical of a National Weather Service warn-  
23 ing. Object matching allows contingency table-based verification statistics to  
24 be used to establish baseline performance metrics for NEWS-e thunderstorm  
25 and mesocyclone forecasts.

26 NEWS-e critical Success Index (CSI) scores of reflectivity (updraft helic-  
27 ity) forecasts decrease from approximately 0.7 (0.4) to 0.4 (0.2) over 3 hours  
28 of forecast time. CSI scores decrease through the forecast period, indicating  
29 that errors have not saturated and skill is retained at 3 hours of forecast time.  
30 Lower verification scores for rotation track forecasts are primarily a result of  
31 a high frequency bias. Comparison of different system configurations used  
32 in 2016 and 2017 show an increase in skill for 2017 reflectivity forecasts, at-  
33 tributable mainly to improvements in the forecast initial condition. A small  
34 decrease in skill in 2017 rotation track forecasts is likely a result of sample  
35 differences between 2016 and 2017. Although large case-to-case variation is  
36 present, evidence is found that NEWS-e forecast skill improves with increas-  
37 ing object size and intensity, as well as in mesoscale environments in which  
38 an enhanced or higher risk of severe thunderstorms was forecast.

## 39 1. Introduction

40 NOAA’s Warn-on-Forecast (WoF) project is tasked with producing probabilistic, short-term  
41 O(0–3 hr) guidance for thunderstorm hazards (Stensrud et al. 2009, 2013). In recent years, pro-  
42 totype WoF systems have demonstrated an ability to produce accurate ensemble forecasts in case  
43 studies of tornado-producing mesocyclones (e.g. Dawson et al. 2012; Yussouf et al. 2013, 2015;  
44 Supinie et al. 2017), severe hail (Snook et al. 2016; Labriola et al. 2017), and flash flooding (Yus-  
45 souf et al. 2016). One system, the NSSL Experimental Warn-on-Forecast System for ensembles  
46 (NEWS-e) has provided ensemble forecasts in real time during the springs of 2016–2017 (Wheat-  
47 ley et al. 2015; Jones et al. 2016). In 2016 and 2017, NEWS-e forecasts were issued up to 17  
48 times daily at 30-minute intervals for a 750 x 750 km domain where severe thunderstorms were  
49 expected. The large amount of forecast data produced during these real time cases makes subjec-  
50 tive verification, which has typically been employed for case studies, difficult and motivates the  
51 development of automated verification techniques for WoF guidance.

52 Automating verification of WoF guidance for thunderstorm hazards presents several challenges.  
53 Firstly, forecasts are issued at convection-allowing scales, typically with  $\sim 3$  km horizontal grid  
54 spacing, which requires the use of spatial verification methods (e.g. Gilleland et al. 2009, 2010)  
55 to avoid double penalties in point verification metrics associated with small displacement errors  
56 (Wilks 2011). Secondly, WoF is interested in predicting localized, rare events occurring in con-  
57 vective storms. These events occur infrequently compared to quantities typically used in model  
58 verification, such as precipitation, even during widespread severe weather outbreaks (e.g. Yussouf  
59 et al. 2015). Finally, phenomena such as mesocyclones are not fully sampled by conventional  
60 observations, which requires development of verification datasets from imperfect proxies of thun-

61 derstorm hazards (Sobash et al. 2011; Skinner et al. 2016; Sobash et al. 2016a,b; Dawson et al.  
62 2017).

63 Verification techniques based on object identification and matching (e.g. Davis et al. 2006a,b;  
64 Ebert and Gallus 2009) are appealing for overcoming the challenges associated with verification  
65 of WoF guidance. Object-based methods are designed to be applicable to non-continuous and non-  
66 traditional features of interest (Davis et al. 2006a). Additionally, object identification and matching  
67 algorithms are adaptable to a variety of user needs. For example, objects may be matched accord-  
68 ing to user-defined total interest values (Davis et al. 2006a,b) and objects derived from different  
69 input fields can be used in verification provided they are consistently defined to isolate features  
70 of interest (Wolff et al. 2014). Finally, object-based methods provide extensive diagnostic in-  
71 formation about forecast and observed objects, allowing specific error sources in forecasts to be  
72 quantified. These advantages have resulted in extensive use of object-based methods for verifica-  
73 tion of convection-allowing model forecasts. Recent examples include verification of quantitative  
74 precipitation estimates (Gallus 2010; Hitchens et al. 2012; Johnson and Wang 2012; Duda and  
75 Gallus 2013; Johnson and Wang 2013; Johnson et al. 2013; Clark et al. 2014; Schwartz et al.  
76 2017), as well as specific features in radar (Burghardt et al. 2014; Pinto et al. 2015; Cai and Du-  
77 mais 2015; Skinner et al. 2016; Sobash et al. 2016b; Burlingame et al. 2017; Jones et al. 2018),  
78 satellite (Griffin et al. 2017a,b), or damage (Clark et al. 2012, 2013; Stratman and Brewster 2017)  
79 proxies.

80 A final complication specific to verification of WoF guidance is that accurate forecasts are  
81 needed on spatial and temporal scales typical of thunderstorm warning products issued by the  
82 National Weather Service. These small time and space scales limit the utility of local storm re-  
83 ports as a verification dataset (Sobash et al. 2011, 2016a,b) owing to errors in the timing, location,  
84 and reporting frequency of severe weather (e.g. Brooks et al. 2003; Doswell et al. 2005; Trapp



85 et al. 2006). In contrast, proxies from Doppler radar observations can be matched to model out-  
86 put with minimal errors in time and space. Additionally, radar data can be used to verify WoF  
87 forecasts in real time, which can be used to provide forecasters with rapidly-updating measures of  
88 forecast performance. These attributes make radar proxies an attractive option for verification of  
89 short-term forecasts of convective storm hazards (Yussouf et al. 2015; Skinner et al. 2016; Dawson  
90 et al. 2017).

91 This study adapts the object-based mesocyclone verification methodology developed by Skinner  
92 et al. (2016) for application to NEWS-e reflectivity and mesocyclone forecasts during 2016<sup>1</sup> and  
93 2017. Verification statistics from 32 total cases are used to establish a baseline of skill for NEWS-e  
94 forecasts of general and severe thunderstorms. Beyond baseline verification statistics aggregated  
95 across all cases, forecast skill is compared for different cases, forecast initialization times, and  
96 object diagnostic properties in order to quantify system performance for differing storm modes  
97 and mesoscale environments. To the authors' knowledge, this study is the first examination of  
98 the skill of Warn-on-Forecast guidance across many cases spanning a variety of storm modes and  
99 mesoscale environments.

100 An object identification and matching strategy for NEWS-e general and severe thunderstorm  
101 forecasts is presented in section 2. Object-based verification metrics and diagnostic properties for  
102 2016 and 2017 NEWS-e composite reflectivity and rotation track forecasts are presented in section  
103 3, including comparisons between different cases, initialization times, and system configurations.  
104 Conclusions, limitations, and recommendations for future research are provided in section 4.

---

<sup>1</sup>A single case from 23 December 2015 is run using the 2016 system configuration and is considered part of the 2016 dataset.

## 2. Methodology

### *a. Description of the Forecast Dataset*

The NEWS-e is an on-demand, ensemble data assimilation and prediction system nested within the High-Resolution Rapid Refresh Ensemble (HRRRE; Dowell et al. 2016). NEWS-e is comprised of an ensemble of 36 WRF-ARW (Skamarock et al. 2008) members with diverse physical parameterizations (Table 1; Wheatley et al. 2015) run over a 750 x 750 km domain with 3-km horizontal grid spacing (Fig. 1). Analyses are initialized at 1800 UTC daily with initial and boundary conditions provided by the HRRRE (Fig. 2) and domain location determined through collaboration with the Storm Prediction Center or as part of the Hazardous Weather Testbed Spring Forecast Experiment (Kain et al. 2003; Gallo and Coauthors 2017). Following initialization, analyses are produced every 15 minutes via assimilation of satellite column integrated liquid or ice water path (Minnis and Coauthors 2011; Jones and Stensrud 2015; Jones et al. 2016), WSR-88D radar reflectivity and radial velocity data, and surface observations using an Ensemble Kalman Filter (EnKF)<sup>2</sup>. Beginning at 1900 UTC, 18-member forecasts with a duration of 180 (90) minutes are issued at the top (bottom) of each hour until 0300 UTC (Fig. 2).

As both NEWS-e and HRRRE are experimental systems being actively developed, several configuration changes were introduced between 2016 and 2017 (Table 2). Differences can be divided into changes in model configuration, changes in HRRRE initial and boundary conditions, and changes in observation processing and assimilation. Model configuration changes from 2016 to 2017 include an upgrade from WRF-ARW version 3.6.1 to 3.8.1 and changing the microphysical parameterization from Thompson (Thompson et al. 2008) to NSSL 2-moment (Mansell et al.

---

<sup>2</sup>The specific EnKF technique is the ensemble adjustment Kalman filter (Anderson 2001) included in the Data Assimilation Research Testbed (DART; Anderson and Collins 2007; Anderson et al. 2009) software. For simplicity, the more general term EnKF is used for the remainder of this manuscript.

2010), which is expected to better represent storm-scale microphysical processes in supercells (Dawson et al. 2010, 2014). Changes to the HRRRE configuration include an expansion of the forecast domain (2017 version shown in Fig. 1), introduction of EnKF-based hourly assimilation of radar reflectivity data, and changes to the observation localization and posterior inflation methodologies (Ladwig et al. 2018). Additionally, 2017 initial conditions for NEWS-e were taken from a 1-hr, 36-member HRRRE forecast initialized at 1700 UTC that provided each NEWS-e analysis member with a unique initial condition. NEWS-e boundary conditions in 2017 were taken from a 9-member HRRRE forecast issued at 1500 UTC and repeated every 9th NEWS-e member. In 2016, NEWS-e initial and boundary conditions were provided by a 3-hr, 18-member HRRRE forecast initialized at 1500 UTC and identical initial and boundary conditions were used for 18 pairs of NEWS-e members. Ensemble spread across these member pairs was produced through diversity in the physics options (Table 1). Finally, assimilation of ASOS observations was performed for 2017 NEWS-e cases 15 minutes past the top of each hour and the methodology for creating Multi-Radar Multi-Sensor (MRMS; Smith et al. 2016) super observations of radar reflectivity data was changed from nearest neighbor to Cressman (Cressman 1959) interpolation. Additional background and details of NEWS-e system configuration are available in Wheatley et al. (2015) and Jones et al. (2016).

NEWS-e forecasts of composite reflectivity and updraft helicity (Kain et al. 2008) in the 2–5 km and 0–2 km layers above ground level (AGL) are examined by this study. These products were selected to test NEWS-e skill in forecasting all thunderstorms (composite reflectivity) and severe thunderstorms (updraft helicity). Examination of updraft helicity calculated over different vertical layers is used to determine if NEWS-e can accurately identify storms producing low-level mesocyclones, which have been found to be the best proxy for tornado occurrence (Trapp et al. 2005). Rotation tracks are used as a final mesocyclone forecast product and are calculated by

150 aggregating 30 minutes of updraft helicity output centered on each 5-minute NEWS-e forecast  
151 timestep.

## 152 *b. Description of the Verification Dataset*

153 Verification of NEWS-e forecasts require proxies for thunderstorm and mesocyclone occurrence  
154 to be derived from WSR-88D data. These proxies are developed using output from the MRMS  
155 system, which provides composite WSR-88D observations across the Continental United States in  
156 real time (Smith et al. 2016).

157 As composite reflectivity observations are available through MRMS, they are an obvious choice  
158 for verification of NEWS-e composite reflectivity forecasts. However, even though the same field  
159 is available in both the forecast and verification datasets it is not an identical, “apples-to-apples”  
160 comparison. Differences between the simulated and observed composite reflectivity will arise  
161 through the model microphysical parameterization, radar sampling differences, and interpolation  
162 of radar data to the model grid. As a result of these differences, simulated and observed composite  
163 reflectivity are treated as different quantities in determining thresholds used for object identifica-  
164 tion (see Section 2c).

165 The verification dataset for mesocyclone forecasts is developed using rotation tracks derived  
166 from MRMS azimuthal wind shear data (Miller et al. 2013). Specifically, maximum range-  
167 corrected MRMS cyclonic azimuthal wind shear (Smith and Elmore 2004; Newman et al. 2013;  
168 Mahalik et al. 2016) in the 0–2 km and 2–5 km layers AGL is calculated every 5 minutes over  
169 the NEWS-e domain. Following quality control and interpolation to the NEWS-e grid, these az-  
170 imuthal wind shear data are aggregated to produce 30-minute rotation tracks for verification of  
171 NEWS-e updraft helicity-based rotation tracks.

172 A challenge in using azimuthal wind shear rotation tracks as a verification dataset is that spurious  
173 observations for rarely occurring phenomena, such as mesocyclones, can have a large impact on  
174 verification metrics. Therefore, extensive quality control is applied to MRMS azimuthal wind  
175 shear fields to mitigate the impact of erroneous observations. Initial quality control is applied  
176 prior to calculation of azimuthal wind shear, with nonmeteorological returns removed by a neural  
177 net trained using polarimetric data (Lakshmanan et al. 2014). Radial velocity data are dealiased  
178 using a modified method of Jing and Weiner (1993) that incorporates near storm environment  
179 soundings from the RAP model. MRMS azimuthal wind shear is then calculated only where the  
180 quality controlled reflectivity is greater than 20 dBZ and blended onto a grid with  $0.01^\circ$  (2016) or  
181  $0.005^\circ$  (2017) latitude/longitude grid spacing. Interpolation of azimuthal wind shear data to the  
182 NEWS-e grid is performed using a Cressman analysis scheme with a 3-km radius of influence.  
183 To be included in the objective analysis, azimuthal wind shear data must be cyclonic<sup>3</sup> and occur  
184 within 20 km of at least 8 MRMS composite reflectivity observations greater than 45 dBZ. At least  
185 four azimuthal wind shear observations must meet these criteria for the grid point to be retained  
186 in the final analysis. The criteria for being retained in the objective analysis of azimuthal wind  
187 shear field are more strict than past studies (Miller et al. 2013) and have been chosen to minimize  
188 spurious values in the output. Finally, regions less than 5-km or greater than 150-km from the  
189 nearest WSR-88D site are removed to mitigate range-related impacts.

### 190 *c. Object Identification*

191 The methodology for object identification in composite reflectivity or rotation track fields is  
192 adapted from the Method for Object-based Diagnostic Evaluation (MODE) software (Davis et al.

---

<sup>3</sup>NEWS-e has produced qualitatively accurate mesoanticyclone forecasts (Jones and Nixon 2017); however, only mesocyclone forecasts are considered by this study.

2006a,b) available in the Model Evaluation Toolkit provided by the National Center for Atmospheric Research. Thunderstorms and mesocyclones are typically sparse, contiguous maxima in both forecast and observation fields, so simple intensity thresholds are used to define object boundaries. However, defining these thresholds is complicated owing to differences in the forecast and verification fields. For example, values that best discriminate mesocyclones in azimuthal wind shear data will be different from the best discriminators in updraft helicity data. To define intensity thresholds that can consistently identify objects in different fields, we assume that a perfect forecast should produce an identical areal footprint in both forecast and verification fields. This assumption allows percentile thresholds (e.g. Mittermaier and Roberts 2010; Dawson et al. 2017) to be used for object identification.

Percentile thresholds are determined using climatologies of forecast and verification fields (Sobash et al. 2016a). These climatologies are sensitive to changes in system configuration, so separate climatologies are constructed for 2016 and 2017 cases (Fig. 3). Each climatology is constructed by aggregating nonzero gridpoint values greater than the domain-wide 99th percentile from each timestep a NEWS-e forecast or interpolated MRMS field is available. These extreme percentile values are used to match thresholds in the forecast and verification fields. The 99.95th percentile value is chosen as a threshold for rotation track objects, which corresponds to 2–5 km updraft helicity and azimuthal wind shear values between 50 and 65  $\text{m}^2 \text{s}^{-2}$  and 0.0035 and 0.005  $\text{s}^{-1}$ , respectively. These updraft helicity values are similar to intensity thresholds used for mesocyclone identification in prior studies (e.g. Kain et al. 2008; Clark et al. 2012; Dawson et al. 2017).

Despite general similarities, clear differences in the climatologies of updraft helicity and azimuthal wind shear are present between 2016 and 2017 (Fig. 3b, c). These differences are attributable to changes in model configuration and the relatively small sample of cases. As updraft

217 helicity is an integrated product of vertical velocity and vertical vorticity, it is sensitive to changes  
218 in the magnitude or alignment of the two input fields. Comparison of cases run using both Thomp-  
219 son and NSSL 2-moment microphysics has revealed that slightly higher values of updraft helicity  
220 are produced by the NSSL 2-moment scheme (not shown). However, comparison of highest SPC  
221 risk and reported tornadoes between 2016 and 2017 cases (Tables 3, 4) reveals that 2016 cases  
222 more often feature tornadic storms in favorable environments. Given the relatively small sam-  
223 ple of cases available, it is therefore likely that changes in the MRMS climatology are primarily  
224 attributable to variation in storm intensity between the years<sup>4</sup>.

225 A composite reflectivity threshold of 45 dBZ is used for NEWS-e output for both 2016 and 2017  
226 and the MRMS threshold is set according to the corresponding percentile (Fig. 3a). As with rota-  
227 tion track output, variation in the composite reflectivity climatology is apparent between 2016 and  
228 2017. Though the MRMS climatology is slightly lower in 2017 than 2016, most of the differences  
229 between the two years are attributable to changes in NEWS-e configuration. Examination of verti-  
230 cal profiles of simulated reflectivity between cases run with both Thompson and NSSL 2-moment  
231 microphysics reveals that the Thompson scheme produces stronger values of simulated reflectivity  
232 above roughly 7 km (Lappin et al. 2018), resulting in much larger maximum NEWS-e composite  
233 reflectivity values in 2016 than 2017. While these differences are most pronounced for NEWS-e  
234 values above  $\sim 50$  dBZ, the MRMS percentile corresponding to 45 dBZ is similar for both 2016  
235 (99.292%) and 2017 (99.374%).

236 The changes in climatologies from year to year illustrate difficulties in establishing an adaptable  
237 object identification methodology for proxy variables such as composite reflectivity or rotation  
238 tracks. The large number of tunable parameters, from quality control of observations through ob-

---

<sup>4</sup>MRMS azimuthal wind shear were merged to a coarser grid in 2016 than 2017; however, differences attributable to MRMS grid spacing are largely smoothed out during interpolation to the NEWS-e grid.

ject identification and matching, are a limitation of object-based verification techniques. Thresholds used in object identification and matching in this study have been determined through trial and error and have been consistently applied in order to compare between different fields and system configurations. Changes to thresholds used for object identification result in different numerical values of verification metrics, but little qualitative change in comparisons between 2016 and 2017 (see Appendix).

Prior to matching forecast and verification objects, a final series of quality control measures are applied in order to minimize retention of spurious objects (Fig. 4). A size threshold of 100 (144) km<sup>2</sup> is applied to rotation track (composite reflectivity) objects. Additionally, rotation track objects are subjected to a continuity threshold of 15 minutes, which requires tracks to consist of input from at least 3 separate timesteps. Finally, objects with a minimum spatial displacement of less than 10 km are merged into a single object.

#### *d. Object Matching and Verification*

Objects in the forecast and verification fields, as well as their associated diagnostic properties, are extracted using the Scikit-image python library (Van der Walt et al. 2014). Forecast and verification objects are then matched according to a total interest score (Davis et al. 2006a,b) adapted from Skinner et al. (2016) using the centroid and minimum spatial displacement and time displacement between object pairs as inputs:

$$TI = \left[ \frac{\left( \frac{cd_{max} - cd}{cd_{max}} \right) + \left( \frac{md_{max} - md}{md_{max}} \right)}{2} \right] \left( \frac{t_{max} - t}{t_{max}} \right) \quad (1)$$

where  $TI$  is the total interest score,  $cd$  the centroid distance between an object pair,  $md$  the minimum distance between an object pair, and  $t$  the time difference between an object pair. The  $max$  subscript indicates the maximum allowable threshold for object matching and is set to 40 km for



centroid and minimum distance and 25 minutes for time displacement. Total interest scores are calculated for each possible pair of forecast and verification objects, with matched pairs requiring a total interest score greater than 0.2, as in Skinner et al. (2016). In cases where multiple forecast objects are matched to a single verification object, only the forecast object with the highest total interest is retained as a match, other objects are reclassified as unmatched.

Calculation of the total interest for this study uses fewer input properties than are typically used in MODE. This simplification is made possible by the generally sparse and contiguous objects in both forecast and verification fields, which allows representative object matching using a small number of input measures (Schwartz et al. 2017). The mean of the two measures of spatial displacement is used as a single input to the final total interest in order to allow matching of objects that may largely overlap but have centroid displacements greater than the allowable threshold, which often occurs for reflectivity objects associated with mesoscale convective systems. As with object identification thresholds, verification scores are sensitive to the maximum allowable offsets in space and time, but qualitative comparisons between datasets remain similar (see Appendix).

Object matching allows matched object pairs to be classified as “hits”, unmatched forecast objects as “false alarms”, and unmatched verification objects as “misses” (Fig. 4). These classifications allow the contingency table-based probability of detection (POD), false alarm ratio (FAR), frequency bias (BIAS), and critical success index (CSI) to be used to quantify the skill of NEWS-e reflectivity and mesocyclone forecasts. Given that object matching does not produce a quantity analogous to correct negatives in the contingency table, verification metrics are limited to those that consider only hits, misses, and false alarms. Additionally, missed verification objects are calculated as the residual of the number of observed objects and number of matched forecast objects at each timestep. This approach results in infrequent occurrences where observed objects are incorrectly classified owing to forecast objects matched across timesteps.

284 Beyond bulk contingency table verification measures, diagnostic features associated with ob-  
285 jects allow specific forecast errors to be identified (Wolff et al. 2014). Specifically, object area,  
286 maximum intensity, and centroid displacement are used in this study to identify variations in fore-  
287 cast skill for different storm modes and intensities and potential phase and storm motion biases,  
288 respectively.

### 289 **3. Object-Based Verification of NEWS-e Forecasts**

#### 290 *a. Comparison of 2016 and 2017 Composite Reflectivity Forecasts*

291 NEWS-e forecasts were produced for a total of 14 cases during 2016 and 18 cases during 2017  
292 across a variety of geographic locations, storm modes, and storm environments (Tables 3, 4).  
293 Variation in cases between years prevents direct comparison of the impacts of NEWS-e system  
294 configuration changes on forecast skill; however, bulk verification metrics for the two years can be  
295 qualitatively compared. Baselines of NEWS-e composite reflectivity forecast skill for 2016 and  
296 2017 have been produced by aggregating all object hits, misses, and false alarms from each case  
297 and ensemble member, then calculating the POD, FAR, BIAS, and CSI at each available forecast  
298 time (Fig. 5).

299 The ability of rapidly-cycling assimilation of radar and satellite data to accurately initialize  
300 individual thunderstorms is evident in the verification metrics as a high probability of detection  
301 and low false alarm ratio in NEWS-e composite reflectivity forecasts (Fig. 5a, c). NEWS-e POD  
302 20 minutes into the forecast is over 0.7 (0.8) for 2017 (2016), with corresponding false alarm ratios  
303 of approximately 0.4 for both years. This initial skill decreases with increasing forecast time, but  
304 does not level off before the end of the forecast period, indicating that forecast errors do not

305 saturate and skill is retained through 3 hours. The POD remains above the FAR for approximately  
306 75 minutes of forecast time for both 2016 and 2017.

307 Despite the generally skillful composite reflectivity forecasts for both years, clear differences  
308 are apparent between 2016 and 2017 (Fig. 5). A positive bias is present during the first forecast  
309 hour for both years, but is more pronounced in 2016 forecasts. This positive bias in 2016 forecasts  
310 results in a higher POD during the first 30 forecast minutes, but 2017 forecasts have a higher POD  
311 for all following times after biases between the years become similar. Furthermore, 2017 forecasts  
312 have a lower false alarm ratio through the duration of the forecast, which combined with the higher  
313 POD at later forecast times, results in higher CSI scores at all forecast times.

314 Examples of the composite reflectivity object distribution from a single forecast with similar  
315 CSI scores to the 2017 ensemble mean are provided in Figure 6. These “paintball” plots illustrate  
316 the accuracy of a NEWS-e reflectivity forecast with CSI scores similar to the yearly mean, with  
317 most ensemble members correctly predicting the position of thunderstorms within a developing  
318 MCS along the Missouri and Arkansas border. In this example, most of the forecast error is driven  
319 by missed objects along the western extent of the domain in eastern Oklahoma. Although some  
320 ensemble members correctly predict the location of these thunderstorms, many do not, particu-  
321 larly for developing convection during the second hour of the forecast (Fig. 6c, d). Several false  
322 alarm objects are also present, mainly in southern Missouri and southeastern Oklahoma; however,  
323 these false alarm objects occur in only a few ensemble members, resulting in low ensemble mean  
324 false alarm ratios. Finally, phase errors are apparent in the forecast of the MCS along the eastern  
325 Missouri and Arkansas borders, with NEWS-e predictions lagging the observed evolution 2 hours  
326 into the forecast (Fig. 6d). Despite these phase errors, many of the ensemble member objects are  
327 classified as matches owing to minimum and centroid distance displacements lower than the pre-  
328 scribed thresholds. This example was selected to illustrate what a NEWS-e forecast that produces

329 CSI values roughly similar to the 2017 mean *can* look like. Many combinations of POD, FAR,  
330 and BIAS can produce similar CSI values and variation is observed across different cases (Fig. 7),  
331 forecasts within a single case, or within the evolution of a single forecast<sup>5</sup>.

332 Case-to-case variation in skill, as well as differences between 2016 and 2017 NEWS-e com-  
333 posite reflectivity forecasts, are apparent comparing performance diagrams (Roebber 2009) 60-  
334 minutes into the forecast of each available case (Fig. 7). With the exception of one outlier, 2017  
335 cases are more clustered, with ensemble mean CSI and frequency values between roughly 0.3–0.6  
336 and 0.75–1.5, respectively. The one outlier case, 2 May 2017, featured isolated storms that initi-  
337 ated after 0100 UTC, resulting in the fewest forecast and observed reflectivity objects from either  
338 year. In contrast, more case-to-case variation is present in 2016 forecasts, with CSI and BIAS  
339 values of approximately 0.2–0.5 and 0.5–2.0, respectively.

340 Changes in NEWS-e performance for different storm modes and environments is examined by  
341 categorizing each case according to the maximum SPC 1630 UTC Day 1 categorical risk within  
342 the NEWS-e domain and subjectively-determined primary storm mode (Tables 3, 4). There is  
343 evidence of stratification of composite reflectivity CSI scores by SPC categorical risk in 2017  
344 forecasts, where an Enhanced risk or higher was present for 6 of the highest 9 scoring cases and  
345 a Slight risk or lower for 7 of the lowest 9 scoring cases. Similar stratification is not apparent for  
346 2016 cases, though the distribution is heavily weighted towards cases with Enhanced risk or higher.  
347 No clear differences in skill are apparent between cases classified as supercellular or mixed/linear  
348 storm mode in either 2016 or 2017.

349 Temporal variation in NEWS-e composite reflectivity forecasts is examined by aggregating ob-  
350 jects across cases for each hourly forecast initialization time (Fig. 8). A decrease in BIAS and

---

<sup>5</sup>At the time of writing, NEWS-e forecast graphics and verification statistics from each case are archived at [www.nssl.noaa.gov/projects/wof/news-e/images.php](http://www.nssl.noaa.gov/projects/wof/news-e/images.php).

351 FAR with increasing forecast initialization time is evident in both 2016 and 2017 cases. These  
352 decreases are coupled with a decrease in POD at later initialization times; however, this decrease  
353 is smaller than decreases in FAR, resulting in a net increase in CSI. These changes with forecast  
354 initialization time likely arise in part from repeated cycles of data assimilation producing more  
355 accurate analyses of existing thunderstorms and suppressing spurious convection. Initiation of  
356 additional convection, which will take several data assimilation cycles to be accurately analyzed  
357 by NEWS-e (e.g. Yussouf and Stensrud 2010) likely contributes to the decrease in POD with  
358 increasing forecast initialization time.

359 Though variation between 2016 and 2017 verification metrics is present for all different initial-  
360 ization times, the largest differences are for forecasts initialized at 2000 and 2100 UTC (Fig. 8).  
361 The CSI of 2017 forecasts at these times is notably higher, at times greater than 0.1, than 2016  
362 forecasts. We surmise that this improvement is likely primarily attributable to upgrades in the  
363 HRRRE between 2016 and 2017, which include hourly ensemble assimilation of radar reflectivity  
364 observations and alterations to the observation localization and posterior inflation methodologies  
365 (Ladwig et al. 2018). These improvements provide NEWS-e forecasts with an improved storm  
366 and mesoscale initial condition that translates to improved NEWS-e performance for early fore-  
367 cast periods.

368 Beyond changes in skill during earlier forecasts, 2016 composite reflectivity forecasts generally  
369 have a higher frequency bias than 2017 forecasts, particularly early in the forecast period (Fig. 8).  
370 This positive bias is additionally evident in bulk (Fig. 5) and case-to-case (Fig. 7) comparisons of  
371 2016 and 2017 forecasts and is primarily a function of different microphysical parameterizations  
372 utilized in 2016 and 2017 (Table 2). The sensitivity of frequency bias to microphysical parameter-  
373 ization is demonstrated by reproducing six cases from 2017 with an identical configuration except

374 that Thompson<sup>6</sup> microphysics is used in place of the NSSL 2-Moment scheme. Cases re-run with  
375 Thompson microphysics all exhibit higher frequency biases in 60-minute composite reflectivity  
376 forecasts than those run with NSSL 2-Moment (Fig. 9). Despite the consistent differences in  
377 frequency bias, compensating variations in POD and success ratio (1 - FAR) occur between the  
378 two sets of experiments, resulting in small and variable changes to the CSI. Composite reflectivity  
379 objects are identified in Thompson experiments using the 2016 NEWS-e reflectivity climatology  
380 (Fig. 3). Since only cases from 2017 were compared, biases will be impacted by differences in  
381 the observed reflectivity climatology between 2016 and 2017. However, the increase in frequency  
382 bias for Thompson runs is exacerbated if either the 2016 or 2017 climatology is applied to both  
383 sets of experiments (not shown) and the results match subjective member-by-member comparisons  
384 between the two sets of experiments, providing confidence that the two schemes produce differing  
385 biases of thunderstorm coverage.

#### 386 *b. Comparison of 2016 and 2017 Updraft Helicity Forecasts*

387 In general, object-based verification scores are lower for mesocyclone forecasts than reflectivity  
388 forecasts (Figs. 10, 11). The CSI for NEWS-e 2–5 km updraft helicity rotation track forecasts  
389 decreases from approximately 0.35–0.45 to 0.2 over the course of a 3 hour forecast during both  
390 2016 and 2017, a reduction of about 0.1 from CSI scores for composite reflectivity forecasts  
391 (Fig. 5). This reduction in CSI is primarily driven by a higher FAR in updraft helicity forecasts  
392 and corresponds to a small positive frequency bias at all forecast times. The positive frequency  
393 bias and increased FAR for 2–5 km rotation track objects indicate that NEWS-e overpredicted  
394 midlevel mesocyclone development in thunderstorms in both 2016 and 2017, especially given

---

<sup>6</sup>An updated, aerosol aware version of the Thompson scheme (Thompson and Eidhammer 2014) was used for these experiments, which is different than the version used for 2016 cases. The impact of changes within the Thompson scheme on NEWS-e forecasts is not known and beyond the scope of this paper.

395 nearly unbiased reflectivity forecasts following the first forecast hour (Fig. 5). Despite generally  
396 lower scores than reflectivity forecasts, the CSI of rotation track forecasts decreases through the  
397 entirety of the 3-hour forecast, suggesting that forecast errors do not saturate and NEWS-e retains  
398 skill through the period.

399 Verification scores for NEWS-e 0–2 km updraft helicity forecasts are generally similar, although  
400 slightly lower, than scores for 2–5 km updraft helicity forecasts (Fig. 11). The number of rota-  
401 tion track objects in the 0–2 km layer is about 5% (20%) lower in 2017 (2016), resulting in a  
402 smaller overprediction bias and reduced POD and FAR. Though fewer low-level rotation track  
403 objects are identified, the strong similarities in verification scores suggest that NEWS-e forecasts  
404 are generally not discriminating between low- and midlevel mesocyclone development. This lack  
405 of discrimination is consistent with prior studies that have found that horizontal grid spacing of 1  
406 km or less is needed to resolve storm-scale processes responsible for low-level mesocyclogenesis  
407 (e.g. Potvin and Flora 2015).

408 NEWS-e 2–5 km updraft helicity forecasts performed slightly better in 2016 than 2017 during  
409 the first hour of the forecast (Fig. 10), exhibiting both a higher probability of detection and lower  
410 false alarm ratio. However, there is large case-to-case variability in forecast performance at 60-  
411 minutes for both 2016 and 2017 (Fig. 12), with CSI and BIAS values ranging from less than 0.1  
412 to greater than 0.4 and roughly 0.25 to greater than 4.0, respectively. Though consistent variation  
413 of forecast skill in different storm modes and environments is not apparent, there is evidence that  
414 supercell cases in a favorable environment most reliably produce high verification scores. For  
415 example, 11 of the 14 highest-scoring cases across both years (roughly equivalent to cases with  
416 CSI greater than 0.35) are from days with Enhanced or greater risk. Furthermore, 5 of the 8  
417 supercell days with enhanced risk or greater are among the 14 highest-scoring cases, including  
418 cases that produced significant tornadoes on 9 May 2016, 24 May 2016, 16 May 2017, and 18

419 May 2017. The higher proportion of supercell cases in a favorable environment in 2016 than 2017  
420 indicates that sample differences<sup>7</sup> may contribute to the improved performance in the first hour of  
421 2016 forecasts.

422 Comparison of 2–5 km rotation track forecast verification metrics from the 6 cases reproduced  
423 using Thompson microphysics (Fig. 13) further suggests variation in skill between 2016 and 2017  
424 forecasts is attributable to sampling differences. Changing the microphysical parameterization  
425 results in small, inconsistent changes to POD, FAR, BIAS, and CSI across the 6 cases. Further-  
426 more, using the 2016 climatological threshold for object identification results in poor scores and  
427 large positive biases greater than 2.0 for all 6 cases, regardless of microphysical parameterization.  
428 This reduction in skill using the 2016 climatology confirms that changes in the updraft helicity  
429 climatology between 2016 and 2017 are primarily driven by differences in the observations, as  
430 opposed to changes in the composite reflectivity climatology, which are primarily driven by the  
431 microphysical parameterization (Figs. 3, 9).

432 In addition to case-to-case variation in verification scores for rotation track forecasts, some cases  
433 exhibit large differences between performance of composite reflectivity and 2–5 km updraft he-  
434 licity forecasts (cf. Figs. 7, 12). In these cases NEWS-e produces generally accurate predictions  
435 of composite reflectivity objects, but less skillful predictions of rotation tracks. Many cases with  
436 the largest reductions in CSI (greater than 0.2) in updraft helicity forecasts are characterized by  
437 predominantly mixed-mode or linear convection, and include 31 March 2016, 3 May 2017, 11  
438 May 2017, 17 May 2017, and 23 May 2017. The reduced performance in rotation track fore-  
439 casts in these cases is typically attributable to either underforecasts of mesocyclones embedded in

---

<sup>7</sup>Though large samples of individual forecast objects are available, many of these objects will be highly correlated owing to the ensemble and high frequency forecast output in NEWS-e. Therefore, sample diversity is better represented by the number of different cases rather than the total number of objects.



mesoscale convective systems or overforecasts of mesocyclones in cellular convection. Examples of the two error sources are provided in Fig. 14, where, despite accurate composite reflectivity forecasts, most ensemble members miss mesocyclone development within an MCS in Iowa (Fig. 14b) or dramatically overpredict mesocyclone development within mixed-mode storms in Texas (Fig. 14d).

Similarly to composite reflectivity forecasts, cycled data assimilation results in a reduction of the BIAS and FAR, and increase in CSI with later forecast initialization times in mesocyclone forecasts (Fig. 15). Differences between 2016 and 2017 are inconsistent and at times highly variable across successive forecasts. However, it appears that 2017 CSI is improved in the 2000 and 2100 UTC forecasts, though to a lesser extent than composite reflectivity forecasts. Additionally, CSI scores for 2016 are higher during the first 30–90 minutes of each forecast from 2200 UTC onward, indicating the improved skill in the first hour of bulk comparisons (Fig. 10) is consistent across most initialization times. Finally, 2016 forecasts initialized at 0200 UTC perform much better than 2017 forecasts. This improvement is not present in 0200 UTC reflectivity forecasts (Fig. 8) and the reasons for the improvement are not clear. However, 4 of the 14 cases from 2016 did not issue forecasts at 0200 UTC (Table 3), which results in far fewer rotation track objects in 2016 than 2017 and will amplify sampling differences between the years.

### *c. Comparison of Object Diagnostic Properties between 2016 and 2017*

Variation of NEWS-e performance with storm characteristics is examined by comparing differences between the size and maximum intensity of matched and false alarm forecast objects. Differences between these diagnostic properties are visualized using scatterplots of composite reflectivity and rotation track objects aggregated from 60-minute NEWS-e forecasts (Fig. 16). Kernel density estimation (KDE) is then used similarly to Anderson-Frey et al. (2016) to highlight

regions within the size and maximum intensity parameter space where object properties occur most often. The KDE technique implemented here applies a Gaussian kernel with a smoothing bandwidth determined from a general optimization algorithm (Scott 1992) to each point within the parameter space. Kernels for each point are summed to provide a measure of the density of points and quantify differences between the distribution of false alarm and matched objects.

Comparison of the size and maximum intensity of NEWS-e reflectivity objects reveal that larger and more intense objects were more likely to be matched to observations in both 2016 and 2017 forecasts (Fig. 16a, b). This result is unsurprising as larger thunderstorms will be better resolved by the 3-km grid spacing employed by NEWS-e and more Doppler radar and satellite observations will be available for assimilation, likely resulting in a more accurate ensemble analysis. In addition to differences between the size and intensity of matched and false alarm objects, differences between the object characteristics in 2016 and 2017 are apparent. As in the model climatologies (Fig. 3), much higher maximum composite reflectivity values are produced by the Thompson microphysical parameterization, with the strongest storms exhibiting values between 70 and 76 dBZ, compared to 58–64 dBZ in NSSL 2-Moment forecasts. Additionally, a small secondary peak in the 2017 object maximum intensity distribution is apparent at roughly 46 dBZ (Fig. 16b). This peak is produced by misidentified objects within the stratiform region of mesoscale convective systems. These spurious objects represent less than 5% of the total number of reflectivity objects in 60-minute forecasts and will minimally impact verification scores, but their presence in NSSL 2-Moment forecasts provides another example of challenges in identifying appropriate thresholds for object-based comparison of different system configurations.

Similarly to reflectivity objects, larger and more intense rotation track objects were more likely to be matched in 2016 forecasts (Fig. 16c), but smaller differences between the distribution of matched and false alarm objects are present in 2017 forecasts (Fig. 16d). However, if 2017

cases are split according to the subjectively defined primary storm mode (Table 4), supercell cases behave similarly to 2016 forecasts, with larger and more intense objects being more likely to be matched to observations (Fig. 16e). In addition, the ratio of matched to false alarm objects in supercell cases from 2017 is higher than for mixed or linear storm modes and similar to the ratio from 2016 cases. This apparent dependence of performance on storm mode provides further evidence that the increased skill during the first hour of updraft helicity forecasts during 2016 is a product of sampling differences between the years rather than changes in model configuration.

Finally, centroid displacement in matched objects is examined to identify potential positive storm motion biases, which have been noted in previous prototype WoF forecasts (Yussouf et al. 2013; Wheatley et al. 2015; Yussouf et al. 2015; Skinner et al. 2016). In contrast with prior studies that found consistent, positive biases in storm speed for forecasts of discrete supercells, large variation in the centroid displacement of matched objects is present in 30-minute NEWS-e forecasts of composite reflectivity and updraft helicity (Fig. 17). Much of this variation results from inclusion of several cases with varying storm modes and coverage. Despite the larger total variation in centroid displacement, north and eastward biases in centroid displacement, consistent with a positive bias in storm speed, are present in 2016 reflectivity forecasts and updraft helicity forecasts from both 2016 and 2017. Though this apparent storm motion bias is consistent with past results and subjective assessment of NEWS-e forecasts, centroid displacement biases can also arise through differences in simulated storm structure (Potvin et al. 2018). For example, changes to reflectivity or rotation track object size with different physical parameterizations will induce changes to object centroid positions and displacement from an observed object. As variation in the distribution of object sizes is noted between 2016 and 2017 for both reflectivity and rotation track objects (Fig. 14), it is unclear to what extent biases in centroid displacement are attributable to errors in storm motion or storm and rotation track structure.

## 4. Conclusions and Future Work

An object-based strategy for verifying Warn-on-Forecast guidance has been presented and applied to 32 cases from 2016 and 2017. Composite reflectivity and updraft helicity-based rotation track forecasts from the NSSL Experimental Warn-on-Forecast System for ensembles are verified against corresponding observations in Multi-Radar Multi-Sensor products on time and space scales typical of National Weather Service warnings. Forecast and verification objects are classified as matched pairs, false alarms, and misses (Fig. 4) allowing contingency table-based metrics to be used to establish a baseline of WoF performance for general and severe thunderstorms. Bulk verification scores from NEWS-e forecasts support the following conclusions:

- Percentile thresholds derived from model climatologies provide a method for prescribing appropriate object identification thresholds to different forecast and verification fields; for example, rotation tracks derived from predicted updraft helicity and observed azimuthal wind shear (Fig. 3).
- Cycled assimilation of Doppler radar and satellite cloud liquid water path observations every 15 minutes will accurately initialize individual thunderstorms within the NEWS-e domain, resulting in POD values greater than 0.7 and FAR values below 0.4 in NEWS-e 30-minute forecasts of composite reflectivity (Fig. 5).
- Critical Success Index scores of NEWS-e composite reflectivity and updraft helicity forecasts decrease through the entirety of 3 hours of forecast time, indicating that forecast errors do not saturate and some skill is retained through the forecast period (Figs. 5, 10).
- NEWS-e composite reflectivity forecasts are more accurate than updraft helicity forecasts, with CSI scores  $\sim 0.1$  higher throughout the forecast period. This reduced performance in

updraft helicity forecasts is primarily a result of overforecasting mesocyclone occurrence (Fig. 10)

- Little difference in NEWS-e forecast skill is evident when considering updraft helicity in the 0–2 km or 2–5 km vertical layers (Figs. 10, 11), indicating that NEWS-e horizontal grid spacing is too coarse to resolve storm-scale processes responsible for development of low-level mesocyclones.

Additionally, the following differences are observed between varying system configurations, storm modes, and storm environments:

- Improvement in composite reflectivity forecasts was noted from 2016 to 2017 and primarily driven by a lower FAR (Fig. 5). The improved performance is attributable to upgrades to the HRRRE, which provides a more accurate initial condition to NEWS-e and results in more accurate early forecasts (Fig. 8) and to implementing the NSSL 2-Moment microphysical parameterization, which reduces a positive frequency bias during the first hour of forecasts (Fig. 9).
- Updraft helicity forecasts during 2016 are more accurate than in 2017 during the first hour of forecast time, with a higher POD and lower FAR (Fig. 10). Inconsistent changes in CSI for 2017 cases rerun with Thompson microphysics (Fig. 13), and a similar skill to 2016 forecasts in 2017 cases with a primarily cellular storm mode (Fig. 16), suggest that improvements in 2016 forecasts are driven by sampling differences between the two years.
- There is tentative evidence that NEWS-e forecasts perform better for larger, more intense storms in favorable environments. The majority of composite reflectivity and updraft helicity cases with the highest CSI contain regions within an Enhanced or higher severe thunderstorm

555 risk in the 1630 UTC SPC Day 1 Convective Outlook (Figs. 7, 12). Additionally, larger  
556 and more intense reflectivity and rotation track objects are more likely to be matched to  
557 observations (Fig. 16).

558 This study has demonstrated the utility of object-based verification for providing a bulk assess-  
559 ment of skill in Warn-on-Forecast guidance, comparing performance across different cases and  
560 system configurations, and providing information on specific forecast errors through examination  
561 of object diagnostic properties. However, there are many limitations to the object-based approach  
562 for short-term, ensemble forecasts of thunderstorm hazards. Object-based verification is highly  
563 customizable, with user-defined thresholds required for object identification and matching (Davis  
564 et al. 2006a). While this flexibility permits application of object-based verification to a wide va-  
565 riety of forecast problems, care must be taken to ensure that appropriate thresholds are used for  
566 consistent object identification and matching in different datasets, particularly for verification of  
567 rare events where small differences in the number of objects identified can dramatically alter ver-  
568 ification scores (Fig. 4). A second limitation to the object-based verification strategy employed  
569 here is that it only provides measures of skill for deterministic forecasts. While this is useful in  
570 establishing general baselines of skill for NEWS-e forecasts, it ignores a fundamental aspect of  
571 Warn-on-Forecast, that guidance should include a measure of uncertainty (Stensrud et al. 2009).  
572 Future work will incorporate additional metrics such as the Brier Skill Score and reliability dia-  
573 grams (Wilks 2011) in order to evaluate probabilistic NEWS-e guidance.

574 The primary limitation of object-based verification specific to this study is limited sample di-  
575 versity across a relatively small number of available cases. Though large numbers of objects are  
576 identified, the ensemble and high frequency nature of NEWS-e forecasts results in strong cor-  
577 relation across forecast objects and variation in model and observation climatologies complicate

578 comparisons between 2016 and 2017 forecasts (Fig. 3). We expect that more regular generation  
579 of real time NEWS-e guidance, as is planned in 2018, will provide a larger sample size of cases  
580 and allow the baseline verification metrics presented here to be refined. Additionally, expanded  
581 computational resources will allow NEWS-e configuration testing across a large sample of prior  
582 cases, permitting hypothesis testing of forecast skill.

583 A final note is that while object-based verification of thunderstorm guidance can provide useful  
584 bulk measures of forecast skill, it does not discriminate between the intensity of different thun-  
585 derstorms. For example, a marginally severe supercell producing a weak rotation track object will  
586 influence verification scores as much as an object associated with a violent tornado. Given the  
587 large numbers of thunderstorms typically present within the NEWS-e domain (e.g. Figs. 6, 14),  
588 changes in forecast quality for the most significant storms for a given case may be masked by  
589 changes to storms that produce limited impacts on life and property. Therefore, subjective verifi-  
590 cation remains indispensable for assessment of forecast skill in case studies of individual storms.

591 *Acknowledgments.* This research was funded by NOAA’s Warn-on-Forecast project with addi-  
592 tional funding provided by the VORTEX-SE project through grant NA16OAR4320115. Partial  
593 funding for this research was also provided by NOAA/Office of Oceanic and Atmospheric Re-  
594 search under NOAA–University of Oklahoma Cooperative Agreement NA11OAR4320072, under  
595 the U.S. Department of Commerce. This paper benefitted from a thoughtful review by Dr. Burkely  
596 Gallo and Drs. Adam Clark, Corey Potvin, and Nusrat Yussouf provided many helpful conversa-  
597 tions over the course of this research. Dr. Darrel Kingfield and Karen Cooper are thanked for  
598 their assistance with MRMS processing. All analyses and visualizations were produced using the  
599 freely-provided Anaconda Python distribution and SciPy, Matplotlib, basemap, netcdf4, sharppy,  
600 scikit-image, and scikit-learn libraries.

## Verification Score Sensitivity to Object Identification and Matching Thresholds

The highly configurable nature of object-based verification measures results in sensitivities of skill scores to user-defined thresholds. The impact of varying the user-defined intensity threshold for object identification and distance threshold for object matching is examined in Figs. A1 and A2.

Variation of the intensity threshold for object identification does result in differences in the probability of detection and false alarm ratios, including changes in comparisons between scores for 2016 and 2017 forecasts (Fig. A1). However, the relative score changes between 2016 and 2017 are attributable to changes in the frequency bias, which produce contrasting changes in POD and FAR that generally result in little net change to the critical success index. An exception is applying the 2016 intensity threshold to 2017 forecasts (Fig. A1 g–i). Using a lower value of updraft helicity for object identification results in approximately 60,000 more rotation track objects in 2017 forecasts that are predominately false alarms, lowering the CSI scores throughout the forecast period. This sensitivity illustrates the importance of considering model climatologies to define representative object identification thresholds when comparing forecast systems with different configurations. Small changes to the percentile threshold produces little relative variation in skill scores between 2016 and 2017 (Fig. A1 j–l) and composite reflectivity forecasts are relatively insensitive to changes in the object identification threshold (not shown), likely owing to small differences between the 2016 and 2017 climatologies below  $\sim 50$  dBZ (Fig. 3a).

As would be expected, increasing the distance threshold for object matching results in corresponding decreases to the false alarm ratio and increases to the probability of detection and critical success index, particularly during the latter portions of the forecast period (Fig. A2). However,



624 there is little relative change between 2016 and 2017 forecasts in any verification metric for either  
625 composite reflectivity or rotation track forecasts.

## 626 **References**

627 Anderson, J. L., 2001: An ensemble adjustment filter for data assimilation. *Mon. Wea. Rev.*, **129**,  
628 2884–2903.

629 Anderson, J. L., and N. Collins, 2007: Scalable implementations of ensemble filter algorithms for  
630 data assimilation. *J. Atmos. Oceanic Technol.*, **59**, 1452–1463.

631 Anderson, J. L., T. Hoar, K. Raeder, H. Liu, N. Collins, R. Torn, and A. Avellano, 2009: The Data  
632 Assimilation Research Testbed: A community facility. *Bull. Amer. Meteor. Soc.*, **90**, 1283–1296.

633 Anderson-Frey, A. K., Y. P. Richardson, A. R. Dean, R. L. Thompson, and B. T. Smith, 2016:  
634 Investigation of near-storm environments for tornado events and warnings. *Wea. Forecasting*,  
635 **31**, 1771–1790.

636 Brooks, H. E., C. A. Doswell III, and M. P. Kay, 2003: Climatological estimates of local daily  
637 tornado probability for the United States. *Wea. Forecasting*, **18**, 626–640.

638 Burghardt, B. J., C. Evans, and P. J. Roebber, 2014: Assessing the predictability of convection  
639 initiation in the High Plains using an object-based approach. *Wea. Forecasting*, **29**, 403–418.

640 Burlingame, B. M., C. Evans, and P. J. Roebber, 2017: The influence of PBL parameterization on  
641 the practical predictability of convection initiation during the Mesoscale Predictability Experi-  
642 ment (MPLEX). *Wea. Forecasting*, **32**, 1161–1183.

643 Cai, H., and R. E. Dumaïs, 2015: Object-based evaluation of a numerical weather prediction  
644 model's performance through storm characteristic analysis. *Wea. and Forecasting*, **31**, 1451–  
645 1468.

646 Clark, A. J., R. G. Bullock, T. L. Jensen, M. Xue, and F. Kong, 2014: Application of object-based  
647 time-domain diagnostics for tracking precipitation systems in convection-allowing models. *Wea.*  
648 *Forecasting*, **29**, 517–542.

649 Clark, A. J., J. Gao, P. T. Marsh, T. Smith, J. S. Kain, J. C. Jr., M. Xue, and F. Kong, 2013: Tornado  
650 pathlength forecasts from 2010 to 2011 using ensemble updraft helicity. *Wea. Forecasting*, **28**,  
651 387–407.

652 Clark, A. J., J. S. Kain, P. T. Marsh, J. C. Jr., M. Xue, and F. Kong, 2012: Forecasting tornado  
653 pathlengths using a three-dimensional object identification algorithm applied to convection-  
654 allowing forecasts. *Wea. Forecasting*, **27**, 1090–1113.

655 Cressman, G. P., 1959: An operational objective analysis system. *Mon. Wea. Rev.*, **87**, 367–374.

656 Davis, C. A., B. G. Brown, and R. G. Bullock, 2006a: Object-based verification of precipitation  
657 forecasts. Part I: Methodology and application to mesoscale rain areas. *Mon. Wea. Rev.*, **134**,  
658 1772–1784.

659 Davis, C. A., B. G. Brown, and R. G. Bullock, 2006b: Object-based verification of precipitation  
660 forecasts. Part II: Application to convective rain systems. *Mon. Wea. Rev.*, **134**, 1785–1795.

661 Dawson, D. T., E. R. Mansell, Y. Jung, L. J. Wicker, M. R. Kumjian, and M. Xue, 2014: Low-  
662 level ZDR signatures in supercell forward flanks: The role of size sorting and melting of hail. *J.*  
663 *Atmos. Sci.*, **71**, 276–299.

664 Dawson, D. T., L. J. Wicker, E. R. Mansell, and R. L. Tanamachi, 2012: Impact of the environ-  
665 mental low-level wind profile on ensemble forecasts of the 4 May 2007 Greensburg, Kansas  
666 tornadic storm and associated mesocyclones. *Mon. Wea. Rev.*, **140**, 696–716.

667 Dawson, D. T., M. Xue, J. A. Milbrandt, and M. K. Yau, 2010: Comparison of evaporation and  
668 cold pool development between single-moment and multimoment bulk microphysics schemes  
669 in idealized simulations of tornadic thunderstorms. *Mon. Wea. Rev.*, **138**, 1152–1171.

670 Dawson, L. C., G. S. Romine, R. J. Trapp, and M. E. Baldwin, 2017: Verifying supercellular ro-  
671 tation in a convection-permitting ensemble forecasting system with radar-derived rotation track  
672 data. *Wea. Forecasting*, **32**, 781–795.

673 Doswell, C. A., H. E. Brooks, and M. P. Kay, 2005: Climatological estimates of daily local non-  
674 tornadic severe thunderstorm probability for the United States. *Wea. Forecasting*, **20**, 577–595.

675 Dowell, D., and Coauthors, 2016: Development of a High-Resolution Rapid Refresh Ensemble  
676 (HRRRE) for severe weather forecasting. *Preprints, 28th Conf. Severe Local Storms*, Portland,  
677 OR, Amer. Meteor. Soc., 8B.2.

678 Duda, J. D., and W. A. Gallus, 2013: The impact of large-scale forcing on skill of simulated  
679 convective initiation and upscale evolution with convection-allowing grid spacings in the WRF.  
680 *Wea. Forecasting*, **28**, 994–1018.

681 Ebert, E. E., and W. A. Gallus, 2009: Toward better understanding of the contiguous rain area  
682 (CRA) method for spatial forecast verification. *Wea. Forecasting*, **24**, 1401–1415.

683 Gallo, B. T., and Coauthors, 2017: Breaking new ground in severe weather prediction: The 2015  
684 NOAA/Hazardous Weather Testbed spring forecasting experiment. *Wea. Forecasting*, **32**, 1541–  
685 1568.

686 Gallus, W. A. J., 2010: Application of object-based verification techniques to ensemble precipita-  
687 tion forecasts. *Wea. Forecasting*, **25**, 144–158.

688 Gilleland, E., D. Ahijevych, B. Brown, and E. Ebert, 2009: Intercomparison of spatial forecast  
689 verification methods. *Wea. Forecasting*, **24**, 1416–1430.

690 Gilleland, E., D. Ahijevych, B. Brown, and E. Ebert, 2010: Verifying forecasts spatially. *Bull.*  
691 *Amer. Meteor. Soc.*, **91**, 1365–1373.

692 Griffin, S. M., J. A. Otkin, C. M. Rozoff, J. M. Sieglaff, L. M. Counce, and C. R. Alexander, 2017a:  
693 Methods for comparing simulated and observed satellite infrared brightness temperatures and  
694 what do they tell us? *Wea. Forecasting*, **32**, 5–25.

695 Griffin, S. M., J. A. Otkin, C. M. Rozoff, J. M. Sieglaff, L. M. Counce, C. R. Alexander, T. L.  
696 Jensen, and J. K. Wolff, 2017b: Seasonal analysis of cloud objects in the High-Resolution  
697 Rapid Refresh (HRRR) model using object-based verification. *J. Appl. Meteor. and Climatology*,  
698 **56**, 2317–2334.

699 Hitchens, N. M., M. E. Baldwin, and R. J. Trapp, 2012: An object-oriented characterization of  
700 extreme precipitation-producing convective systems in the midwestern United States. *Mon. Wea.*  
701 *Rev.*, **140**, 1356–1366.

702 Jing, Z., and G. Weiner, 1993: Two-dimensional dealiasing of Doppler velocities. *J. Atmos.*  
703 *Oceanic Tech.*, **10**, 798–808.

704 Johnson, A., and X. Wang, 2012: Verification and calibration of neighborhood and object-based  
705 probabilistic precipitation forecasts from a multimodel convection-allowing ensemble. *Mon.*  
706 *Wea. Rev.*, **140**, 3054–3077.

707 Johnson, A., and X. Wang, 2013: Object-based evaluation of a storm-scale ensemble during the  
708 2009 NOAA Hazardous Weather Testbed Spring Experiment. *Mon. Wea. Rev.*, **141**, 1079–1098.

709 Johnson, A., X. Wang, F. Kong, and M. Xue, 2013: Object-based evaluation of the impact of  
710 horizontal grid spacing on convection-allowing forecasts. *Mon. Wea. Rev.*, **141**, 3413–3425.

711 Jones, T. A., K. Knopfmeier, D. Wheatley, G. Creager, P. Minnis, and R. Palikondo, 2016:  
712 Storm-scale data assimilation and ensemble forecasting with the NSSL Experimental Warn-  
713 on-Forecast System. Part 2: Combined radar and satellite data experiments. *Wea. Forecasting*,  
714 **31**, 297–327.

715 Jones, T. A., and C. Nixon, 2017: Short-term forecasts of left-moving supercells from an experi-  
716 mental Warn-on-Forecast system. *J. Operational Meteor.*, **5**, 161–170.

717 Jones, T. A., and D. J. Stensrud, 2015: Assimilating cloud water path as a function of model cloud  
718 microphysics in an idealized simulation. *Mon. Wea. Rev.*, **143**, 2052–2081.

719 Jones, T. A., X. Wang, P. S. Skinner, A. Johnson, and Y. Wang, 2018: Assimilation of GOES-13  
720 imager clear-sky water vapor ( $6.5\ \mu\text{m}$ ) radiances into a Warn-on-Forecast system. *Mon. Wea.*  
721 *Rev.*, **145**, In Review.

722 Kain, J. S., P. R. Janish, S. J. Weiss, R. S. Schneider, M. E. Baldwin, and H. E. Brooks, 2003:  
723 Collaboration between forecasters and research scientists at the NSSL and SPC: The spring  
724 program. *Bull. Amer. Meteor. Soc.*, **84**, 1797–1806.

725 Kain, J. S., and Coauthors, 2008: Some practical considerations regarding horizontal resolution in  
726 the first generation of operational convection-allowing NWP. *Wea. Forecasting*, **23**, 931–952.

- 727 Labriola, J., N. Snook, Y. Jung, B. Putnam, and M. Xue, 2017: Ensemble hail prediction for the  
728 storms of 10 May 2010 in south-central Oklahoma using single- and double-moment micro-  
729 physical schemes. *Mon. Wea. Rev.*, **145**, 4911–4936.
- 730 Ladwig, T. T., and Coauthors, 2018: Development of the High-Resolution Rapid Refresh Ensem-  
731 ble HRRRE toward an operational convective-allowing ensemble data assimilation and forecast  
732 system. *Preprints, 6th Symp. on the Weather, Water, and Climate Enterprise*, Austin, TX, Amer.  
733 Meteor. Soc., TJ1.2.
- 734 Lakshmanan, V., C. Karstens, J. Krause, and L. Tang, 2014: Quality control of weather radar data  
735 using polarimetric variables. *J. Atmos. Oceanic Technol.*, **31**, 1234–1249.
- 736 Lappin, F. M., D. M. Wheatley, K. H. Knopfmeier, and P. S. Skinner, 2018: An evaluation  
737 of changes to the NSSL Experimental Warn-on-Forecast System for ensembles (NEWS-e) in  
738 spring 2017. *Preprints, 22nd Conf. on Integrated Observing and Assimilation Systems for the*  
739 *Atmosphere, Oceans, and Land Surface*, Austin, TX, Amer. Meteor. Soc., 170.
- 740 Mahalik, M. C., B. R. Smith, D. M. Kingfield, K. L. Ortega, T. M. Smith, and K. L. Elmore, 2016:  
741 Improving NSSL azimuthal shear calculations using an updated derivation and range-based  
742 corrections. *Preprints, 28th Conf. Severe Local Storms*, Portland, OR, Amer. Meteor. Soc., 182.
- 743 Mansell, E. R., C. L. Ziegler, and E. C. Bruning, 2010: Simulated electrification of a small thun-  
744 derstorm with two-moment bulk microphysics. *J. Atmos. Sci.*, **67**, 171–194.
- 745 Miller, M. L., V. Lakshmanan, and T. M. Smith, 2013: An automated method for depicting meso-  
746 cyclone paths and intensities. *Wea. Forecasting*, **28**, 570–585.

747 Minnis, and Coauthors, 2011: CERES edition-2 cloud property retrievals using TRMM VIRS  
748 and Terra and Aqua MODIS data—Part I: Algorithms. *IEEE Trans. Geosci. Remote Sens.*, **49**,  
749 4374–4400.

750 Mittermaier, M., and N. Roberts, 2010: Intercomparison of spatial forecast methods: Identifying  
751 skillful spatial scales using the fractions skill score. *Wea. Forecasting*, **25**, 343–354.

752 Newman, J. F., V. Lakshmanan, P. L. Heinselman, M. B. Richman, and T. M. Smith, 2013: Range-  
753 correcting azimuthal shear in Doppler radar data. *Wea. Forecasting*, **28**, 194–211.

754 Pinto, J. O., J. A. Grim, and M. Steiner, 2015: Assessment of the High-Resolution Rapid Refresh  
755 model’s ability to predict mesoscale convective systems using object-based evaluation. *Wea.*  
756 *Forecasting.*, **30**, 892–913.

757 Potvin, C. A., J. R. Carley, A. J. Clark, L. J. Wicker, J. S. Kain, A. R. Reinhart, and P. S. Skin-  
758 ner, 2018: Inter-model storm-scale comparisons from the 2017 HWT Spring Forecasting Ex-  
759 periment. *Preprints, Eighth Conf. on Transition of Research to Operations*, Ausin, TX, Amer.  
760 Meteor. Soc., 3B.5.

761 Potvin, C. A., and M. L. Flora, 2015: Sensitivity of idealized supercell simulations to horizontal  
762 grid spacing: Implications for Warn-on-Forecast. *Mon. Wea. Rev.*, **143**, 2998–3024.

763 Roebber, P. J., 2009: Visualizing multiple measures of forecast quality. *Wea. Forecasting.*, **24**,  
764 601–608.

765 Schwartz, C. S., G. S. Romine, K. R. Fossell, R. A. Sobash, and M. L. Weisman, 2017: Toward  
766 1-km ensemble forecasts over large domains. *Mon. Wea. Rev.*, **145**, 2943–2969.

767 Scott, D. W., 1992: *Multivariate Density Estimation: Theory, Practice, and Visualization*. John  
768 Wiley and Sons, 360 pp.

769 Skamarock, W. C., and Coauthors, 2008: A description of the advanced research wrf version 3.  
 770 ATC TN-475+STR, National Center for Atmospheric Research, 113 pp.

771 Skinner, P. S., L. Wicker, D. M. Wheatley, and K. H. Knopfmeier, 2016: Application of two  
 772 spatial verification methods to ensemble forecasts of low-level rotation. *Wea. Forecasting*, **31**,  
 773 713–735.

774 Smith, T. M., and K. L. Elmore, 2004: The use of radial velocity derivatives to diagnose rotation  
 775 and divergence. *Preprints, 11th Conf. on Aviation, Range, and Aerospace*, Hyannis, MA, Amer.  
 776 Meteor. Soc., P5.6.

777 Smith, T. M., and Coauthors, 2016: Multi-Radar Multi-Sensor (MRMS severe weather and avia-  
 778 tion products: Initial operating capabilities. *Bull. Amer. Meteor. Soc.*, **97**, 1617–1630.

779 Snook, N., Y. Jung, J. Brotzge, B. Putnam, and M. Xue, 2016: Prediction and ensemble forecast  
 780 verification of hail in the supercell storms of 20 May 2013. *Wea. Forecasting*, **31**, 811–825.

781 Sobash, R. A., J. S. Kain, D. R. Bright, A. R. Dean, and M. C. Coniglio, 2011: Probabilistic  
 782 forecast guidance for severe thunderstorms based on the identification of extreme phenomena  
 783 in convection-allowing model forecasts. *Wea. Forecasting*, **26**, 714–728.

784 Sobash, R. A., G. S. Romine, C. S. Schwartz, D. J. Gagne, and M. L. Weisman, 2016a: Explicit  
 785 forecasts of low-level rotation from convection-allowing models for next-day tornado predic-  
 786 tion. *Wea. Forecasting*, **31**, 255–271.

787 Sobash, R. A., C. S. Schwartz, G. S. Romine, K. R. Fossell, and M. L. Weisman, 2016b: Severe  
 788 weather prediction using storm surrogates from an ensemble forecasting system. *Wea. Forecast-*  
 789 *ing*, **31**, 255–271.



- 790 Stensrud, D. J., and Coauthors, 2009: Convective-scale Warn-On-Forecast system: A vision for  
791 2020. *Bull. Amer. Meteor. Soc.*, **90**, 1487–1499.
- 792 Stensrud, D. J., and Coauthors, 2013: Progress and challenges with Warn-On-Forecast. *Atmos.*  
793 *Res.*, **123**, 2–16.
- 794 Stratman, D. R., and K. A. Brewster, 2017: Sensitivities of 1-km forecasts of 24 May 2011 tornadic  
795 supercells to microphysics parameterizations. *Mon. Wea. Rev.*, **145**, 2697–2721.
- 796 Supinie, T. A., N. Yussouf, Y. Jung, M. Xue, J. Cheng, and S. Wang, 2017: Comparison of the  
797 analyses and forecasts of a tornadic supercell storm from assimilating phased-array radar and  
798 WSR-88D observations. *Wea. Forecasting*, **32**, 1379–1401.
- 799 Thompson, G., and T. Eidhammer, 2014: A study of aerosol impacts on clouds and precipitation  
800 development in a large winter cyclone. *J. Atmos. Sci.*, **71**, 3636–3658.
- 801 Thompson, G., P. R. Field, R. M. Rasmussen, and W. D. Hall, 2008: Explicit forecasts of winter  
802 precipitation using an improved bulk microphysics scheme. Part II: Implementation of a new  
803 snow parameterization. *Mon. Wea. Rev.*, **136**, 5095–5115.
- 804 Trapp, R. J., G. J. Stumpf, and K. L. Manross, 2005: A reassessment of the percentage of tornadic  
805 mesocyclones. *Wea. Forecasting*, **20**, 680–687.
- 806 Trapp, R. J., D. M. Wheatley, N. T. Atkins, R. W. Przybylinski, and R. Wolf, 2006: Buyer beware:  
807 Some words of caution on the use of severe wind reports in postevent assessment and research.  
808 *Wea. Forecasting*, **21**, 408–415.
- 809 Van der Walt, S., J. L. Schonberger, J. Nunez-Iglesias, F. Boulogne, J. D. Warner, N. Yager,  
810 E. Gouillart, and T. Yu, 2014: Scikit-image: Image processing in Python. *PeerJ*, **2**, e453.

811 Wheatley, D. M., K. H. Knopfmeier, T. A. Jones, and G. J. Creager, 2015: Storm-scale data  
 812 assimilation and ensemble forecasting with the NSSL Experimental Warn-on-Forecast System.  
 813 Part 1: Radar data experiments. *Wea. Forecasting*, **30**, 1795–1817.

814 Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. Elsevier, 627 pp.

815 Wolff, J. K., M. Harrold, T. Fowler, J. H. Gotway, L. Nance, and B. G. Brown, 2014: Beyond  
 816 the basics: Evaluating model-based precipitation forecasts using traditional, spatial, and object-  
 817 based methods. *Wea. and Forecasting*, **29**, 1451–1472.

818 Yussouf, N., D. C. Dowell, L. J. Wicker, K. H. Knopfmeier, and D. M. Wheatley, 2015: Storm-  
 819 scale data assimilation and ensemble forecasts for the 27 April 2011 severe weather outbreak in  
 820 Alabama. *Mon. Wea. Rev.*, **143**, 3044–3066.

821 Yussouf, N., J. S. Kain, and A. J. Clark, 2016: Short-term probabilistic forecasts of the 31 May  
 822 2013 Oklahoma tornado and flash flood event using a continuous-update-cycle storm-scale en-  
 823 semble system. *Wea. Forecasting*, **31**, 957–983.

824 Yussouf, N., E. R. Mansell, L. J. Wicker, D. M. Wheatley, and D. J. Stensrud, 2013: The ensemble  
 825 Kalman filter analyses and forecasts of the 8 May 2003 Oklahoma City tornadic supercell storms  
 826 using single- and double-moment microphysics schemes. *Mon. Wea. Rev.*, **141**, 3388–3412.

827 Yussouf, N., and D. J. Stensrud, 2010: Impact of phased-array radar observations over a short as-  
 828 simulation period: Observing system simulation experiments using an Ensemble Kalman Filter.  
 829 *Mon. Wea. Rev.*, **138**, 517–538.

830 **LIST OF TABLES**

831	<b>Table 1.</b>	Physical parameterization options for NEWS-e forecast members during 2016	
832		and 2017 (adapted from Wheatley et al. 2015, their Table 2). Planetary bound-	
833		ary layer (PBL) options include the Yonsei University (YSU), Mellor–Yamada–	
834		Janjic (MYJ), and Mellor–Yamada–Nakanashi–Niino (MYNN) schemes,	
835		which are paired with either Dudhia and Rapid Radiative Transfer Model	
836		(RRTM) or the Rapid Radiative Transfer Model–Global (RRTMG) parame-	
837		terizations for shortwave and longwave radiation. All members utilize the RAP	
838		land surface parameterization. Physics options for NEWS-e analysis members	
839		19–36 are repeated (e.g. member 19 would have the same options as member	
840		1). . . . .	39
841	<b>Table 2.</b>	Changes in NEWS-e system configuration between 2016 and 2017. Additional	
842		changes to the HRRRE configuration are discussed in Section 2a. . . . .	40
843	<b>Table 3.</b>	Summary of 2016 NEWS-e cases. For each date the available forecast period,	
844		satellite data availability, maximum Storm Prediction Center (SPC) risk from	
845		the 1630 outlook within the NEWS-e domain, number of SPC archived tor-	
846		nado reports within the domain and forecast period, primary states effected,	
847		and predominant storm mode are provided. . . . .	41
848	<b>Table 4.</b>	Same as Table 3, except for 2017 cases. Asterisks indicate cases reproduced	
849		using Thompson microphysics. . . . .	42

850 TABLE 1. Physical parameterization options for NEWS-e forecast members during 2016 and 2017 (adapted  
851 from Wheatley et al. 2015, their Table 2). Planetary boundary layer (PBL) options include the Yonsei University  
852 (YSU), Mellor–Yamada–Janjic (MYJ), and Mellor–Yamada–Nakanashi–Niino (MYNN) schemes, which are  
853 paired with either Dudhia and Rapid Radiative Transfer Model (RRTM) or the Rapid Radiative Transfer Model–  
854 Global (RRTMG) parameterizations for shortwave and longwave radiation. All members utilize the RAP land  
855 surface parameterization. Physics options for NEWS-e analysis members 19–36 are repeated (e.g. member 19  
856 would have the same options as member 1).

Member	PBL	Shortwave Radiation	Longwave Radiation
1	YSU	Dudhia	RRTM
2	YSU	RRTMG	RRTMG
3	MYJ	Dudhia	RRTM
4	MYJ	RRTMG	RRTMG
5	MYNN	Dudhia	RRTM
6	MYNN	RRTMG	RRTMG
7	YSU	Dudhia	RRTM
8	YSU	RRTMG	RRTMG
9	MYJ	Dudhia	RRTM
10	MYJ	RRTMG	RRTMG
11	MYNN	Dudhia	RRTM
12	MYNN	RRTMG	RRTMG
13	YSU	Dudhia	RRTM
14	YSU	RRTMG	RRTMG
15	MYJ	Dudhia	RRTM
16	MYJ	RRTMG	RRTMG
17	MYNN	Dudhia	RRTM
18	MYNN	RRTMG	RRTMG

857      TABLE 2. Changes in NEWS-e system configuration between 2016 and 2017. Additional changes to the  
858 HRRRE configuration are discussed in Section 2a.

	2016	2017
WRF-ARW Version	3.6.1	3.8.1
Microphysics	Thompson	NSSL 2-Moment
Initial Conditions	3-hr HRRRE 1500 UTC Forecast (18 members)	1-hr HRRRE 1700 UTC Forecast (36 members)
Boundary Conditions	HRRRE 1500 UTC Forecast (18 members)	HRRRE 1500 UTC Forecast (9 members)
ASOS Assimilation	No	Hourly
Reflectivity Super Observations	Nearest Neighbor Interpolation	Cressman Objective Analysis

859 TABLE 3. Summary of 2016 NEWS-e cases. For each date the available forecast period, satellite data avail-  
860 ability, maximum Storm Prediction Center (SPC) risk from the 1630 outlook within the NEWS-e domain,  
861 number of SPC archived tornado reports within the domain and forecast period, primary states effected, and  
862 predominant storm mode are provided.

Date	Forecast Period	Satellite DA	SPC Outlook	Tornado Reports	Primary States Affected	Primary Storm Mode
23 December 2015	1900–0100 UTC	No	Moderate	24	AL, MS, TN	Supercell
31 March 2016	1900–0130 UTC	No	Enhanced	24	AL, MS, TN	Mixed
10 April 2016	1900–0300 UTC	No	Enhanced	0	OK, TX	Linear
29 April 2016	1900–2330 UTC	No	Slight	0	AL, MS	Linear
07 May 2016	1900–0300 UTC	Yes	Slight	15	CO, KS	Mixed
08 May 2016	1900–0300 UTC	Yes	Enhanced	9	KS, OK	Supercell
09 May 2016	1900–0100 UTC	Yes	Enhanced	16	AR, KS, OK	Supercell
10 May 2016	1900–0300 UTC	Yes	Enhanced	19	IL, IN, KY	Mixed
16 May 2016	1900–0300 UTC	Yes	Enhanced	10	OK, TX	Linear
17 May 2016	1900–0300 UTC	Yes	Enhanced	1	TX	Mixed
22 May 2016	1900–0300 UTC	Yes	Enhanced	38	KS, OK, TX	Supercell
23 May 2016	1900–0300 UTC	Yes	Enhanced	5	OK, TX	Supercell
24 May 2016	1900–0300 UTC	Yes	Enhanced	29	CO, KS, NE, OK	Supercell
25 May 2016	1900–0300 UTC	Yes	Slight	14	KS, OK	Supercell

863

864

TABLE 4. Same as Table 3, except for 2017 cases. Asterisks indicate cases reproduced using Thompson microphysics.

Date	Forecast Period	Satellite DA	SPC Outlook	Tornado Reports	Primary States Affected	Primary Storm Mode
01 May 2017	1900–0300 UTC	Yes	Enhanced	6	NY, PA	Linear
02 May 2017	1900–0300 UTC	Yes	Slight	0	OK, TX	Supercell
03 May 2017	1900–0300 UTC	Yes	Enhanced	2	LA, TX	Linear
04 May 2017	1900–0300 UTC	Yes	Marginal	11	GA, SC	Mixed
08 May 2017	1900–0300 UTC	Yes	Slight	1	CO, NM	Supercell
09 May 2017*	1900–0300 UTC	Yes	Slight	6	NM, TX	Supercell
11 May 2017	1900–0300 UTC	Yes	Enhanced	11	AR, LA, OK, TX	Mixed
15 May 2017	1900–0300 UTC	Yes	Slight	0	CO, KS, NE	Mixed
16 May 2017*	1900–0300 UTC	Yes	Moderate	26	KS, OK, TX	Supercell
17 May 2017*	1900–0300 UTC	Yes	Enhanced	17	IA, IL, MN, WI	Mixed
18 May 2017*	1900–0300 UTC	Yes	High	34	KS, OK, TX	Supercell
19 May 2017	1900–0300 UTC	Yes	Enhanced	4	OK, TX	Mixed
22 May 2017	1900–0300 UTC	Yes	Slight	0	NM, TX	Supercell
23 May 2017*	1900–0300 UTC	Yes	Slight	0	TX	Mixed
25 May 2017	1900–0300 UTC	Yes	Slight	2	CO, KS	Supercell
26 May 2017	1900–0300 UTC	Yes	Slight	8	CO, KS	Supercell
27 May 2017*	1900–0300 UTC	Yes	Moderate	8	AR, MO, OK	Mixed
30 May 2017	1900–0300 UTC	Yes	Slight	1	MD, PA, VA	Mixed

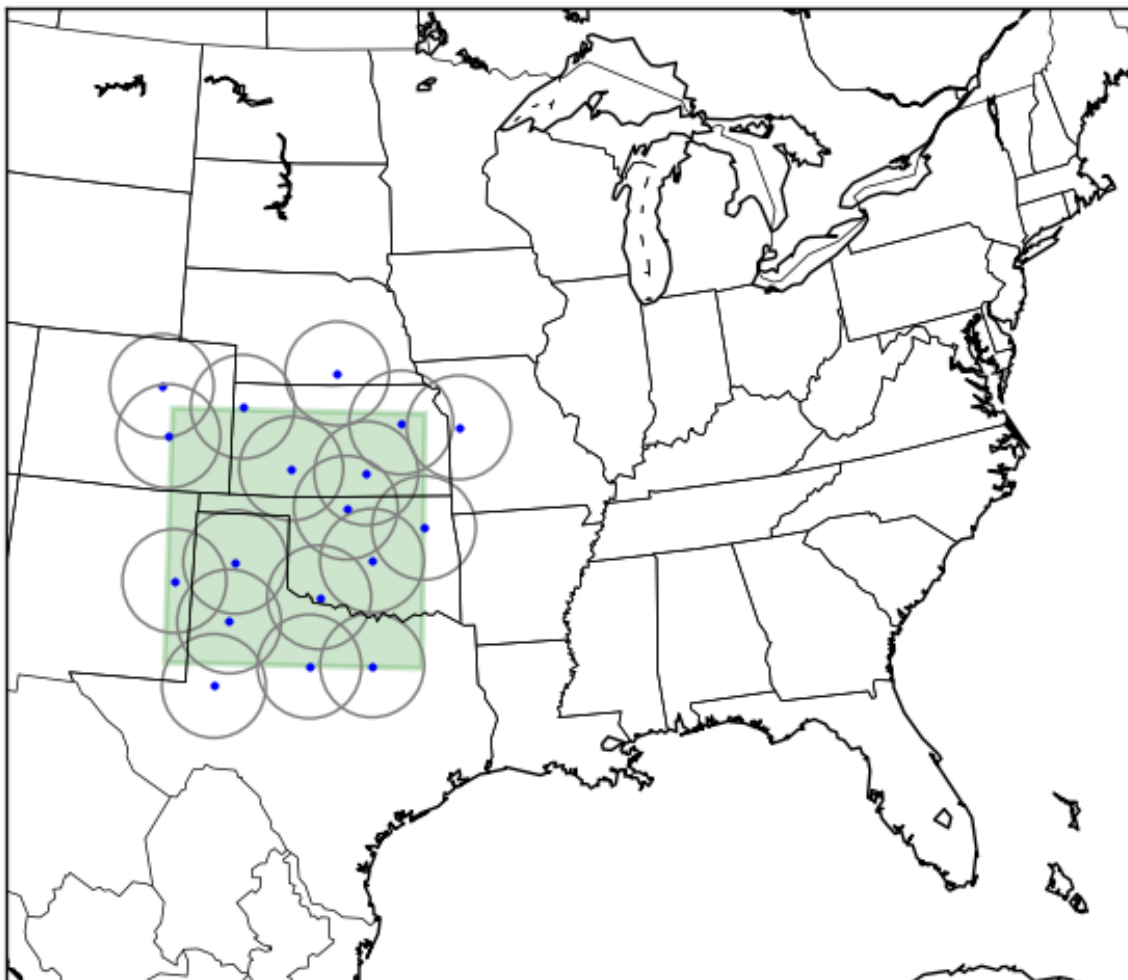
## LIST OF FIGURES

<b>Fig. 1.</b>	Example NEWS-e domain from 16 May 2017. The map shown corresponds to the HRRRE domain, with the nested NEWS-e domain shaded green. WSR-88D sites whose data are assimilated into NEWS-e are marked by blue dots with 150-km range rings drawn in gray.	45
<b>Fig. 2.</b>	Schematic of NEWS-e system configuration for 2017.	46
<b>Fig. 3.</b>	Climatologies of forecast and verification datasets for (blue) 2016 cases and (orange) 2017 cases. Scatter plots show the the 99.1st through 99.98th percentile values for (a) composite reflectivity (dBZ), (b) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ) or azimuthal wind shear (AWS; $\text{s}^{-1}$ ), and (c) 0–2 km updraft helicity or azimuthal wind shear. Thresholds used for object identification are marked by horizontal and vertical lines.	47
<b>Fig. 4.</b>	Schematic depicting the object matching and verification process. Initial thresholded fields from the (a) forecast from a single ensemble member and (d) observations are subjected to size and continuity quality control thresholds prior to (b, e) object identification. (c) Forecast objects are matched to verification objects according to prescribed spatiotemporal displacement thresholds with matched pairs being considered, hits, unmatched forecast objects false alarms, and unmatched verification objects misses. This classification of objects allows the (f) standard contingency table metrics probability of detection (POD), false alarm ratio (FAR), frequency bias (BIAS), and critical success index (CSI) to be calculated to quantify forecast skill.	48
<b>Fig. 5.</b>	Time series of object-based (a) POD, (b) BIAS, (c) FAR, and (d) CSI for composite reflectivity forecasts from (blue) 2016 and (orange) 2017. Individual ensemble members are plotted with thin lines and the ensemble mean in bold. Ensemble means are calculated as the mean of verification metrics from each ensemble member. The first and last 20 minutes of the forecast are masked so that only forecast times where objects could be matched in time as well as space are considered. The total number of objects from each year is annotated.	49
<b>Fig. 6.</b>	Paintball plots of composite reflectivity objects (a) 30, (b) 60, (c) 90, and (d) 120 minutes into forecasts initialized at 0100 UTC on 28 May 2017. Colored shading indicates NEWS-e member forecast objects, with different colors assigned to each ensemble member, and dark gray shading observed objects. Regions shaded light gray are less than 5 km or greater than 150 km from the nearest WSR-88D and not considered in verification. Ensemble mean POD, FAR, BIAS, and CSI scores are provided in the upper right of each panel.	50
<b>Fig. 7.</b>	Performance diagrams (Roebber 2009) for 60-minute composite reflectivity forecasts from each case during (a) 2016 and (b) 2017. Small circles indicate scores of individual ensemble members and large circles represent the ensemble mean from each case. Cases are numbered according to the legend provided below each plot and color coded according to maximum SPC risk in the NEWS-e domain and storm mode. The total number of objects identified for each case is provided following each date in the legend.	51
<b>Fig. 8.</b>	Time series of the object-based ensemble mean (a) POD, (b) FAR, (c) BIAS, and (d) CSI for composite reflectivity forecasts aggregated for each forecast initialization hour between 2000 and 0200 UTC. Scores from 2017 (2016) forecasts are plotted in orange (blue) and every other forecast is plotted using lighter, dashed lines in order to improve readability. As in Fig. 5, the first and last 20 minutes of each forecast are masked. The total number of objects for each forecast initialization hour is annotated in panel a.	52



908	<b>Fig. 9.</b>	As in Fig. 7 except for 60-minute composite reflectivity forecasts from (orange) 6 cases in 2017 and (blue) the same 6 cases re-run using Thompson microphysics. The 2016 reflectivity climatology was used to identify objects in the forecasts using Thompson microphysics.	53
912	<b>Fig. 10.</b>	As in Fig. 5 except for 2–5 km updraft helicity forecasts.	54
913	<b>Fig. 11.</b>	As in Fig. 5 except for 0–2 km updraft helicity forecasts.	55
914	<b>Fig. 12.</b>	As in Fig. 7 except for 60-minute 2–5 km updraft helicity forecasts.	56
915	<b>Fig. 13.</b>	As in Fig. 9 except for 2–5 km updraft helicity forecasts. Note that the 2017 2–5 km updraft helicity climatology is used to define rotation track objects in both the Thompson and NSSL 2-Moment experiments.	57
918	<b>Fig. 14.</b>	As in Fig. 6 except for (a, c) composite reflectivity and (b, d) rotation track objects 60-minutes into forecasts initialized at 2300 UTC on (a, b) 17 May 2017 and (c, d) 23 May 2017. POD, FAR, BIAS, and CSI scores for each forecast are provided in the lower left of each panel. Note that some forecast rotation track objects in (d) are matched to observed objects at different times, resulting in a FAR less than 1.0 despite no observed objects being present at the forecast time plotted.	58
924	<b>Fig. 15.</b>	As in Fig. 8 except for 2–5 km updraft helicity forecasts.	59
925	<b>Fig. 16.</b>	Scatterplots of the parameter space of object area and maximum intensity for 60-minute NEWS-e forecasts of (a, b) composite reflectivity (dBZ) and (c–f) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ) during (a, c) 2016, (b, d) 2017, and 2017 cases classified as (e) supercell or (f) mixed/linear mode. Matched objects are plotted in orange and false alarm objects in blue with the total number of objects in each category listed in the lower right. Kernel density estimate contours of the 95th, 97.5th, 99th, and 99.9th percentile values of each distribution are overlain to illustrate differences between matched and false alarm distributions. Every third reflectivity object is plotted to improve clarity.	60
933	<b>Fig. 17.</b>	Scatterplots of the east-west and north-south centroid displacements (km) of matched objects for 30-minute NEWS-e forecasts of (a) composite reflectivity (dBZ) and (b) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ). Objects from 2016 (2017) are plotted in blue (orange) and the total number of objects for each year is listed in the lower left. Kernel density estimate contours are overlain as in Fig. 16 and every third reflectivity object is plotted to improve clarity.	61
938	<b>Fig. A1.</b>	Time series of (a, d, g, j) probability of detection, (b, e, h, k) false alarm ratio, and (c, f, i, l) critical success index for NEWS-e 2-5 km rotation track forecasts. The intensity threshold used to identify forecast and observed rotation track objects is varied between the (a–c) 99.95th percentile from each year’s climatology (same as Fig. 7), (d–f) the 99.95th percentile from the 2017 climatology only, (g–i) the 99.95th percentile from the 2016 climatology only, and (j–l) the 99.9th percentile from each year’s climatology. Individual ensemble member scores are plotted in thin orange (blue) lines with thick orange (blue) lines representing the ensemble mean for 2017 (2016) NEWS-e forecasts.	62
946	<b>Fig. A2.</b>	As in Fig. A1, except for the composite reflectivity objects and the maximum distance threshold for object matching is varied from (a–c) 20 km, (d–f) 30 km, (g–i) 40 km (same as Fig. 5), and (j–l) 60 km.	63

### 3-km HRRRE background and nested NEWS-e grid



Radar locations within NEWS-e grid shown as blue dots with 150-km range rings

949 FIG. 1. Example NEWS-e domain from 16 May 2017. The map shown corresponds to the HRRRE domain,  
950 with the nested NEWS-e domain shaded green. WSR-88D sites whose data are assimilated into NEWS-e are  
951 marked by blue dots with 150-km range rings drawn in gray.

## 2017 Configuration

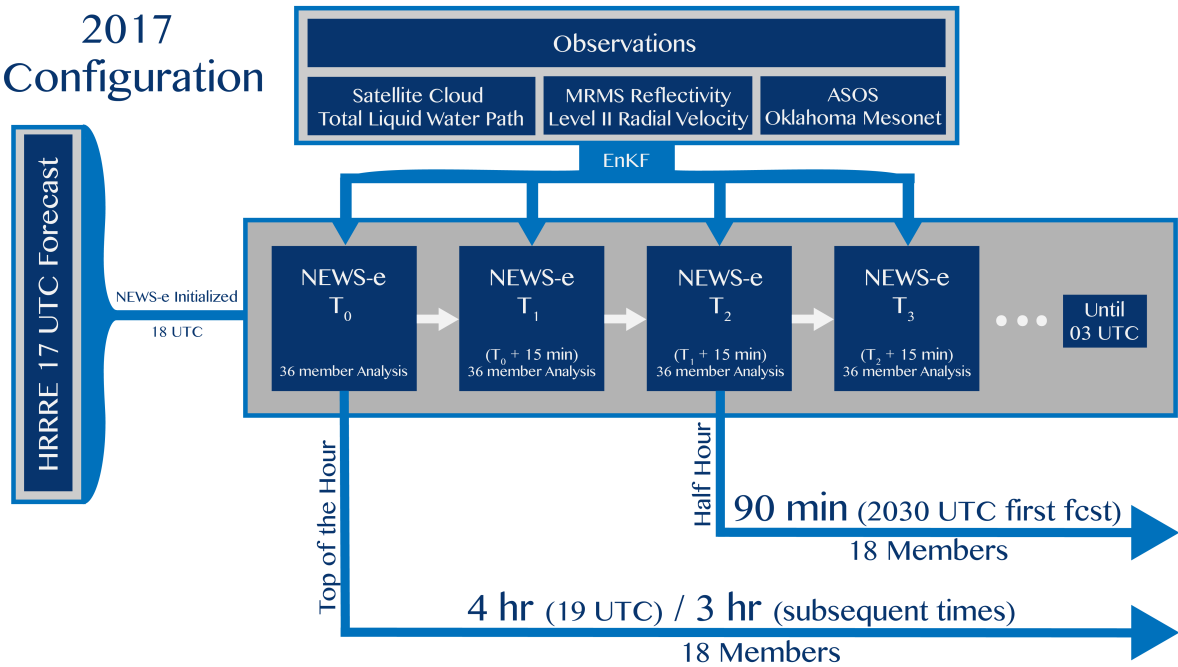


FIG. 2. Schematic of NEWS-e system configuration for 2017.

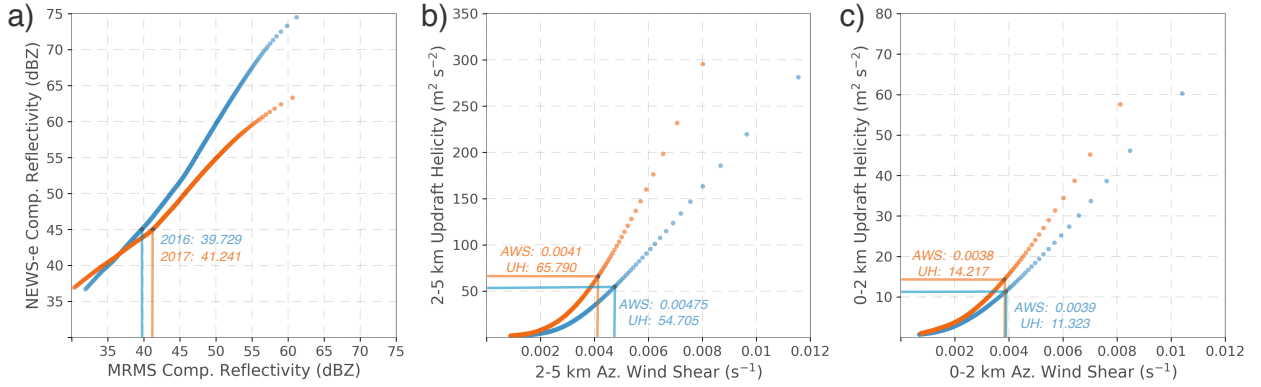


FIG. 3. Climatologies of forecast and verification datasets for (blue) 2016 cases and (orange) 2017 cases. Scatter plots show the the 99.1st through 99.98th percentile values for (a) composite reflectivity (dBZ), (b) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ) or azimuthal wind shear (AWS;  $\text{s}^{-1}$ ), and (c) 0–2 km updraft helicity or azimuthal wind shear. Thresholds used for object identification are marked by horizontal and vertical lines.

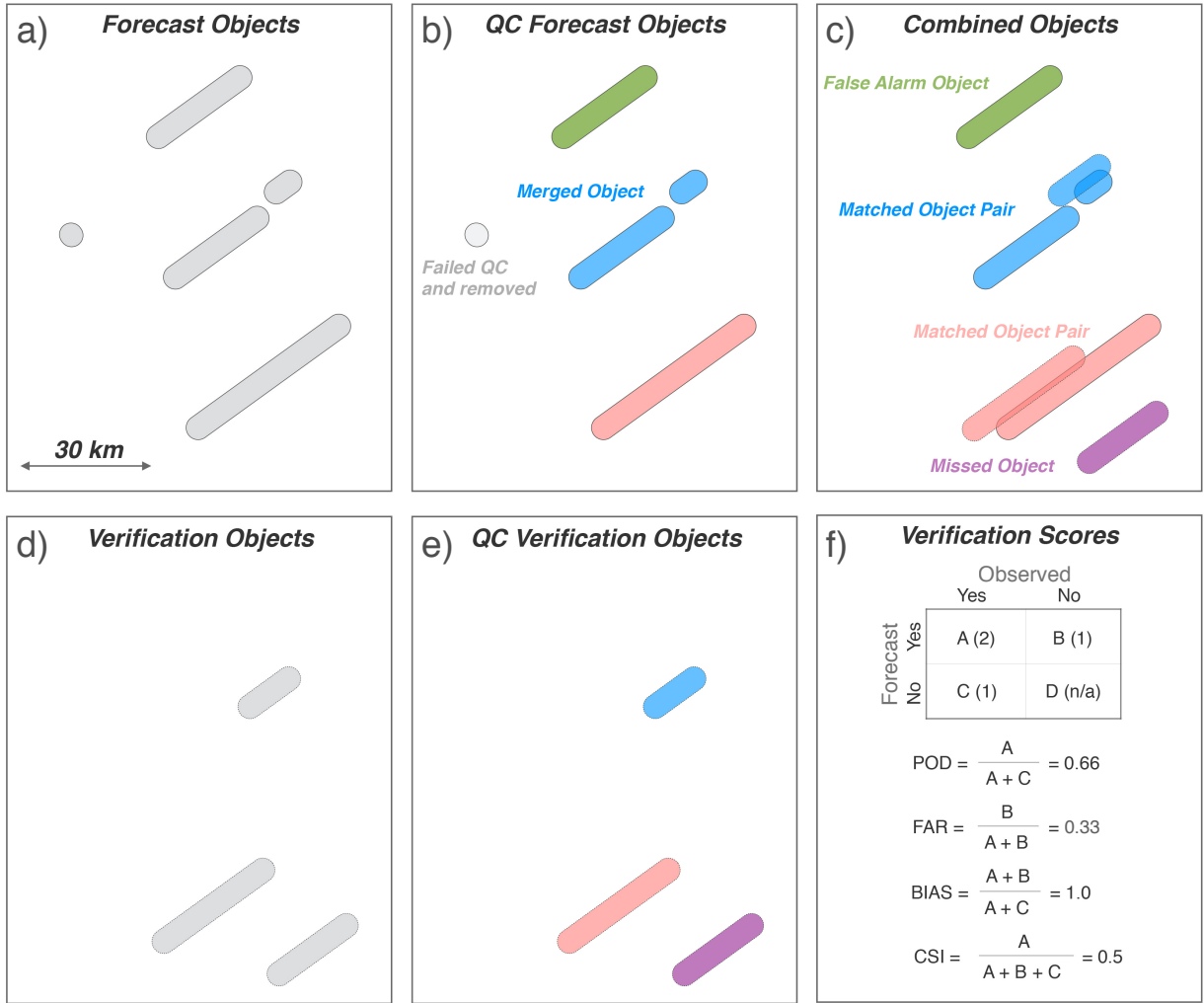


FIG. 4. Schematic depicting the object matching and verification process. Initial thresholded fields from the (a) forecast from a single ensemble member and (d) observations are subjected to size and continuity quality control thresholds prior to (b, e) object identification. (c) Forecast objects are matched to verification objects according to prescribed spatiotemporal displacement thresholds with matched pairs being considered, hits, unmatched forecast objects false alarms, and unmatched verification objects misses. This classification of objects allows the (f) standard contingency table metrics probability of detection (POD), false alarm ratio (FAR), frequency bias (BIAS), and critical success index (CSI) to be calculated to quantify forecast skill.

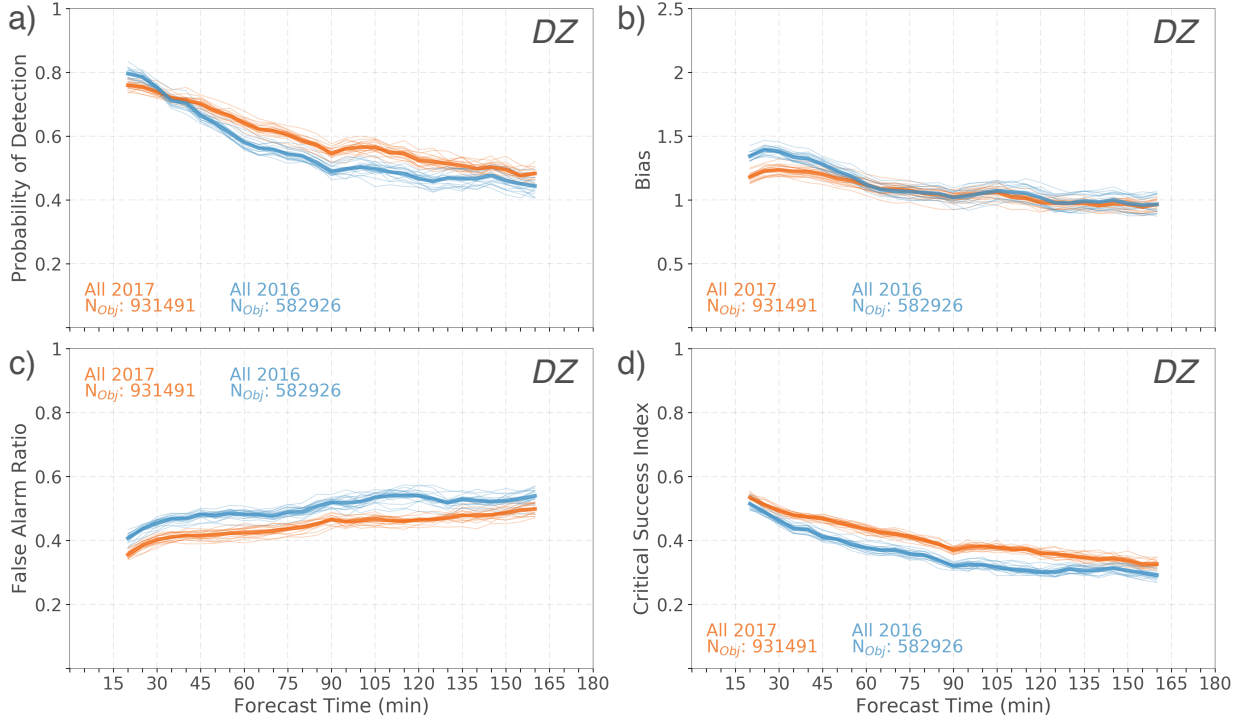


FIG. 5. Time series of object-based (a) POD, (b) BIAS, (c) FAR, and (d) CSI for composite reflectivity forecasts from (blue) 2016 and (orange) 2017. Individual ensemble members are plotted with thin lines and the ensemble mean in bold. Ensemble means are calculated as the mean of verification metrics from each ensemble member. The first and last 20 minutes of the forecast are masked so that only forecast times where objects could be matched in time as well as space are considered. The total number of objects from each year is annotated.

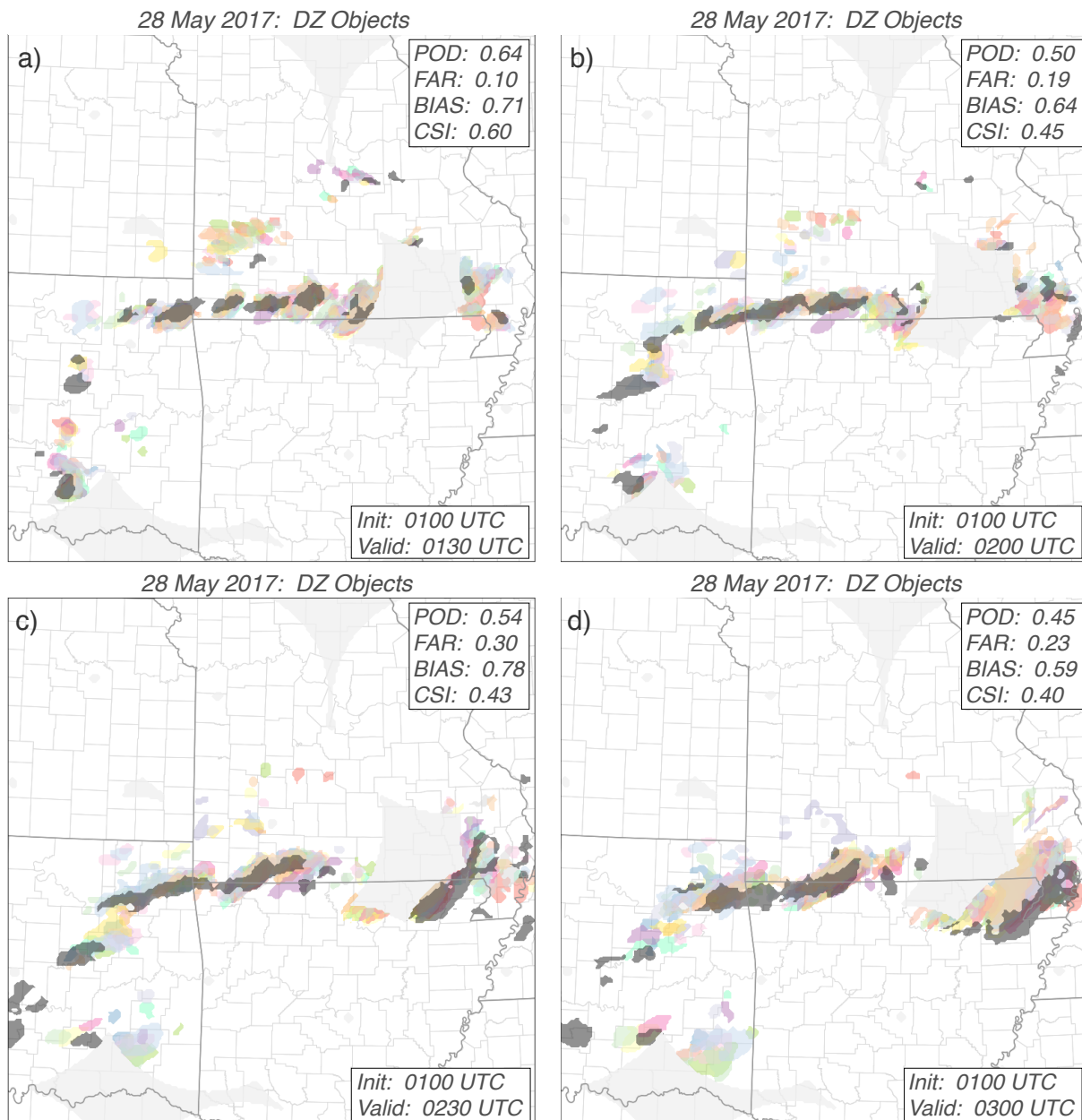


FIG. 6. Paintball plots of composite reflectivity objects (a) 30, (b) 60, (c) 90, and (d) 120 minutes into forecasts initialized at 0100 UTC on 28 May 2017. Colored shading indicates NEWS-e member forecast objects, with different colors assigned to each ensemble member, and dark gray shading observed objects. Regions shaded light gray are less than 5 km or greater than 150 km from the nearest WSR-88D and not considered in verification. Ensemble mean POD, FAR, BIAS, and CSI scores are provided in the upper right of each panel.

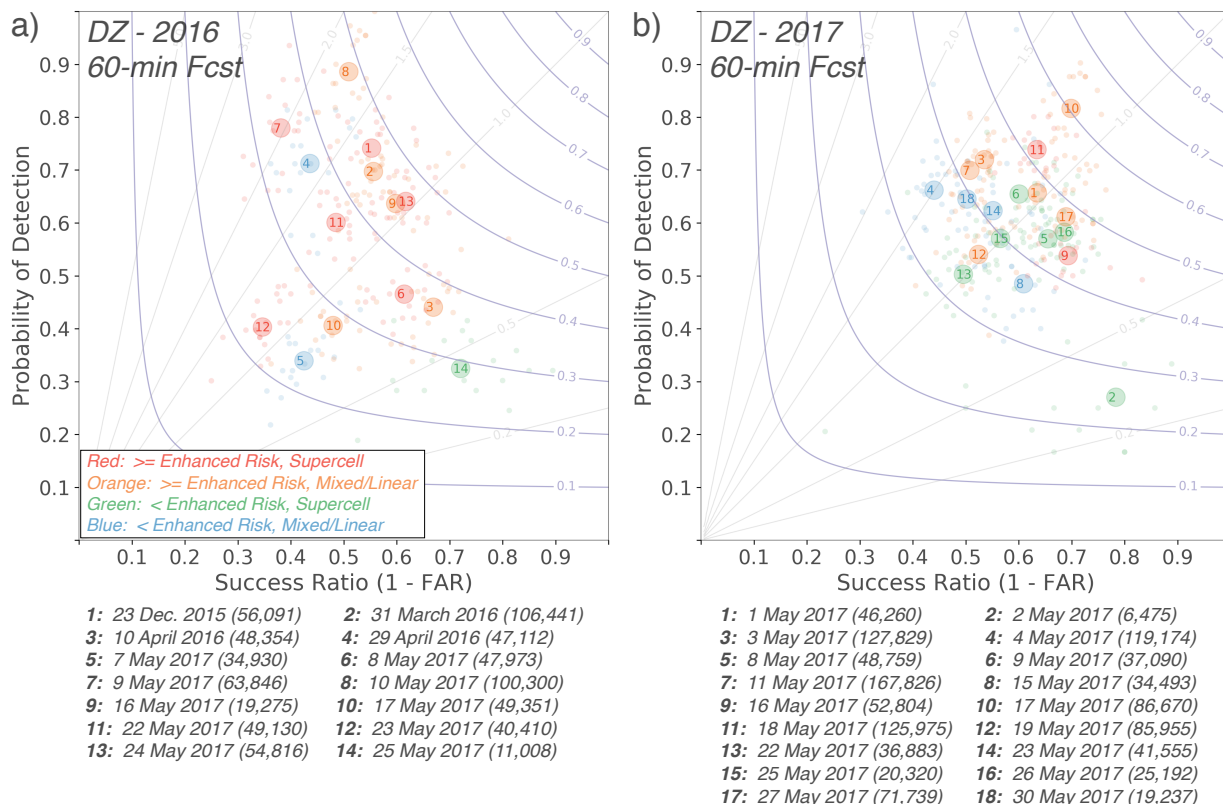


FIG. 7. Performance diagrams (Roebber 2009) for 60-minute composite reflectivity forecasts from each case during (a) 2016 and (b) 2017. Small circles indicate scores of individual ensemble members and large circles represent the ensemble mean from each case. Cases are numbered according to the legend provided below each plot and color coded according to maximum SPC risk in the NEWS-e domain and storm mode. The total number of objects identified for each case is provided following each date in the legend.



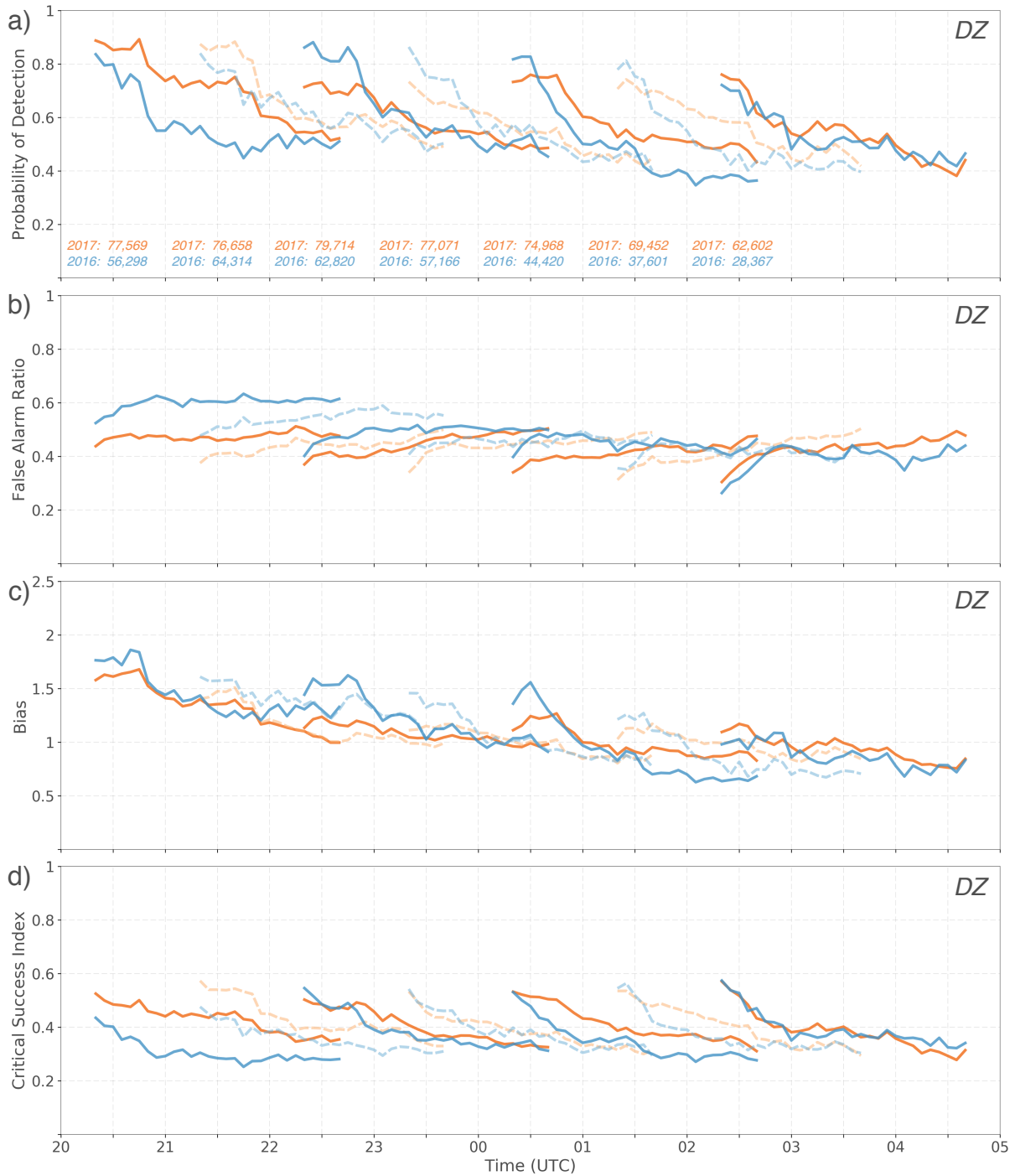


FIG. 8. Time series of the object-based ensemble mean (a) POD, (b) FAR, (c) BIAS, and (d) CSI for composite reflectivity forecasts aggregated for each forecast initialization hour between 2000 and 0200 UTC. Scores from 2017 (2016) forecasts are plotted in orange (blue) and every other forecast is plotted using lighter, dashed lines in order to improve readability. As in Fig. 5, the first and last 20 minutes of each forecast are masked. The total number of objects for each forecast initialization hour is annotated in panel a.

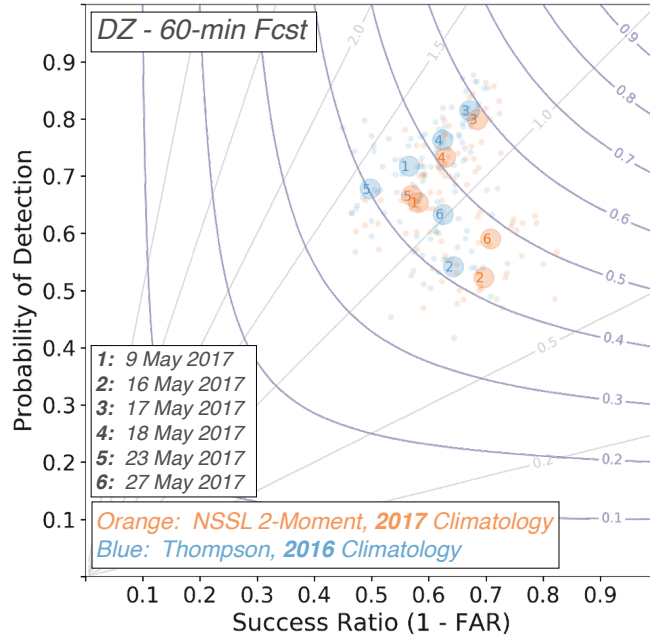


FIG. 9. As in Fig. 7 except for 60-minute composite reflectivity forecasts from (orange) 6 cases in 2017 and (blue) the same 6 cases re-run using Thompson microphysics. The 2016 reflectivity climatology was used to identify objects in the forecasts using Thompson microphysics.

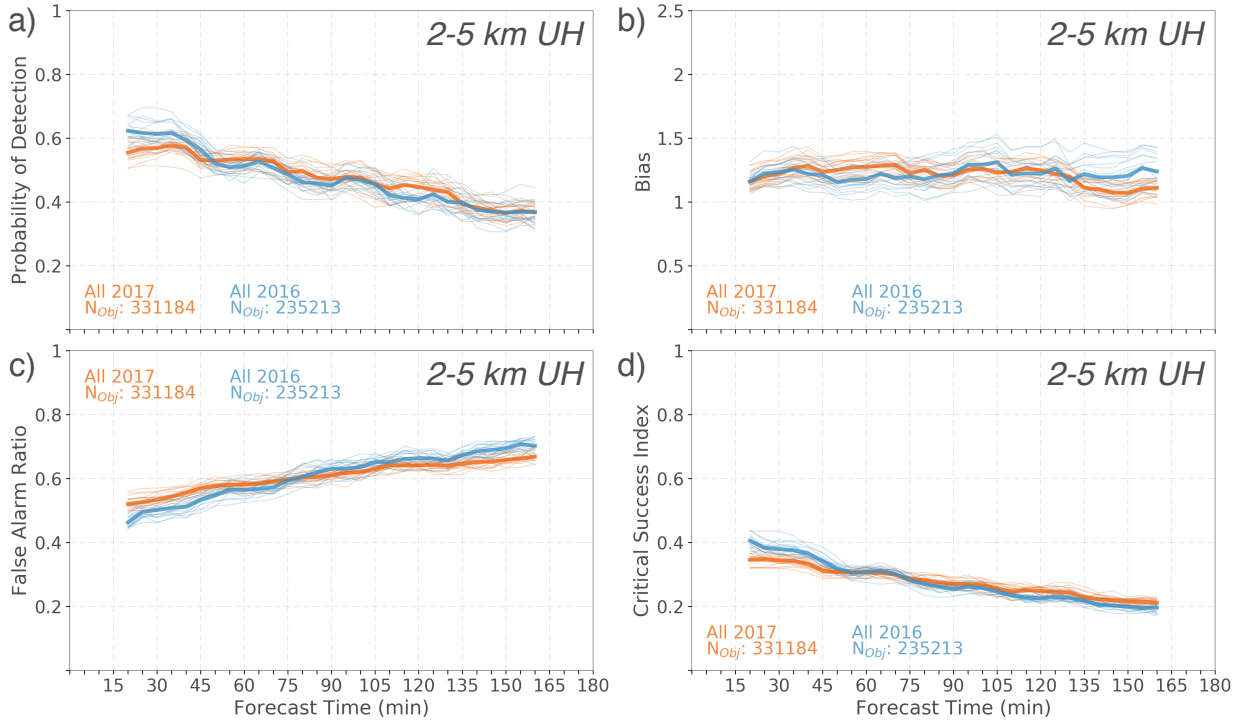


FIG. 10. As in Fig. 5 except for 2–5 km updraft helicity forecasts.

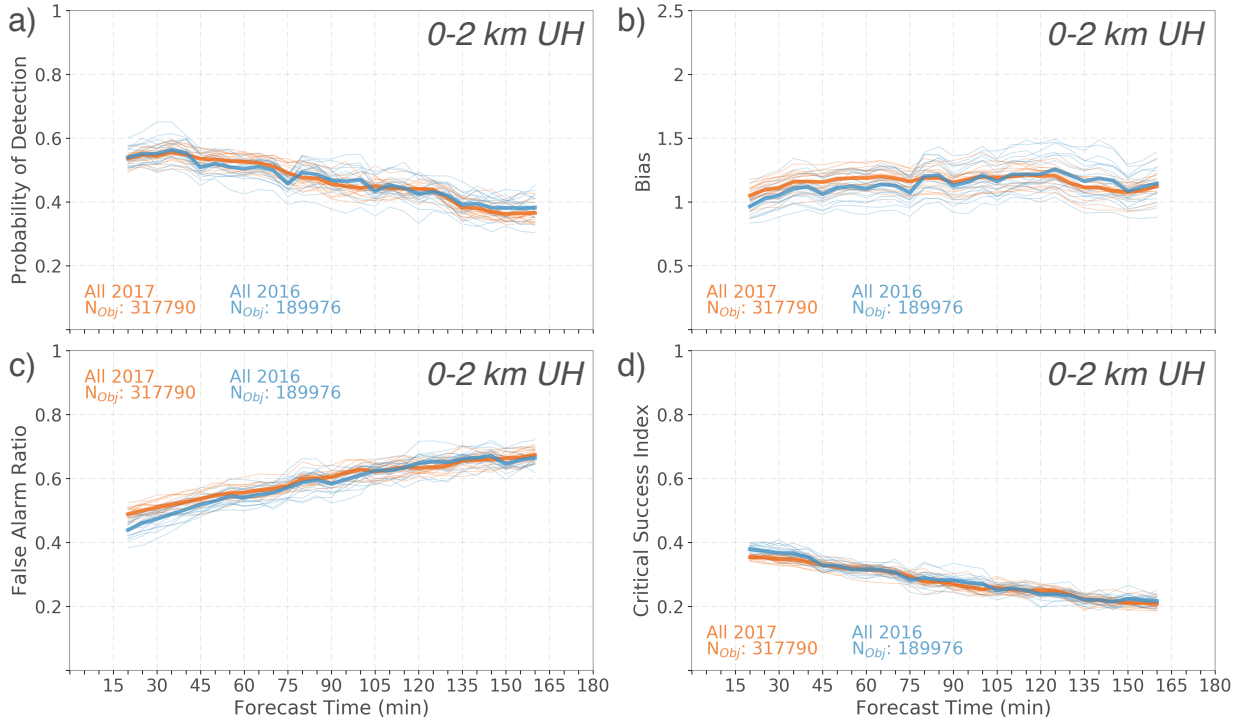


FIG. 11. As in Fig. 5 except for 0–2 km updraft helicity forecasts.

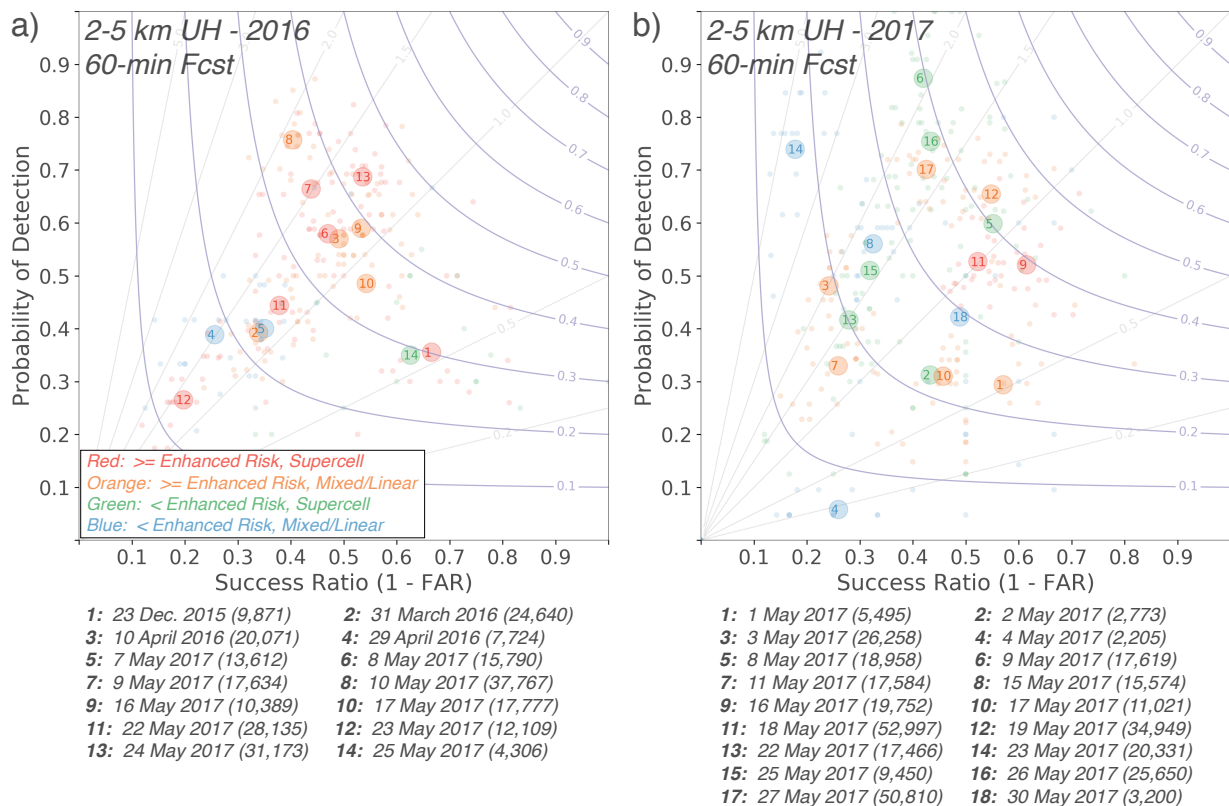


FIG. 12. As in Fig. 7 except for 60-minute 2–5 km updraft helicity forecasts.

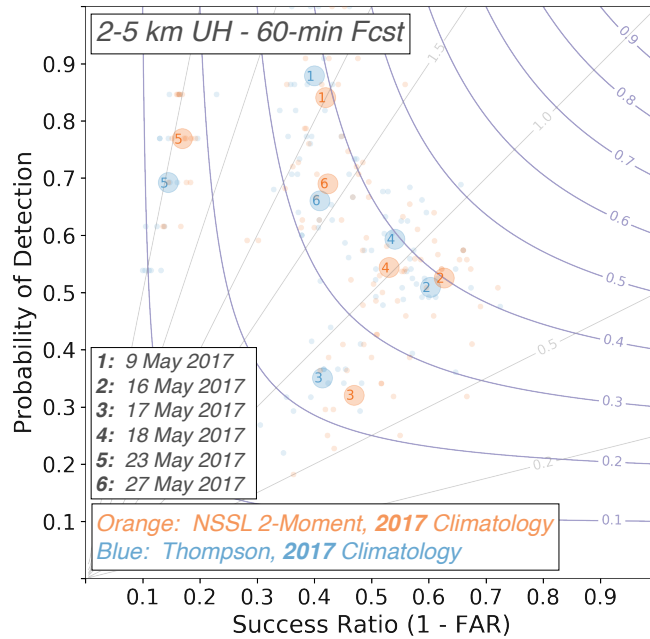


FIG. 13. As in Fig. 9 except for 2–5 km updraft helicity forecasts. Note that the 2017 2–5 km updraft helicity climatology is used to define rotation track objects in both the Thompson and NSSL 2-Moment experiments.

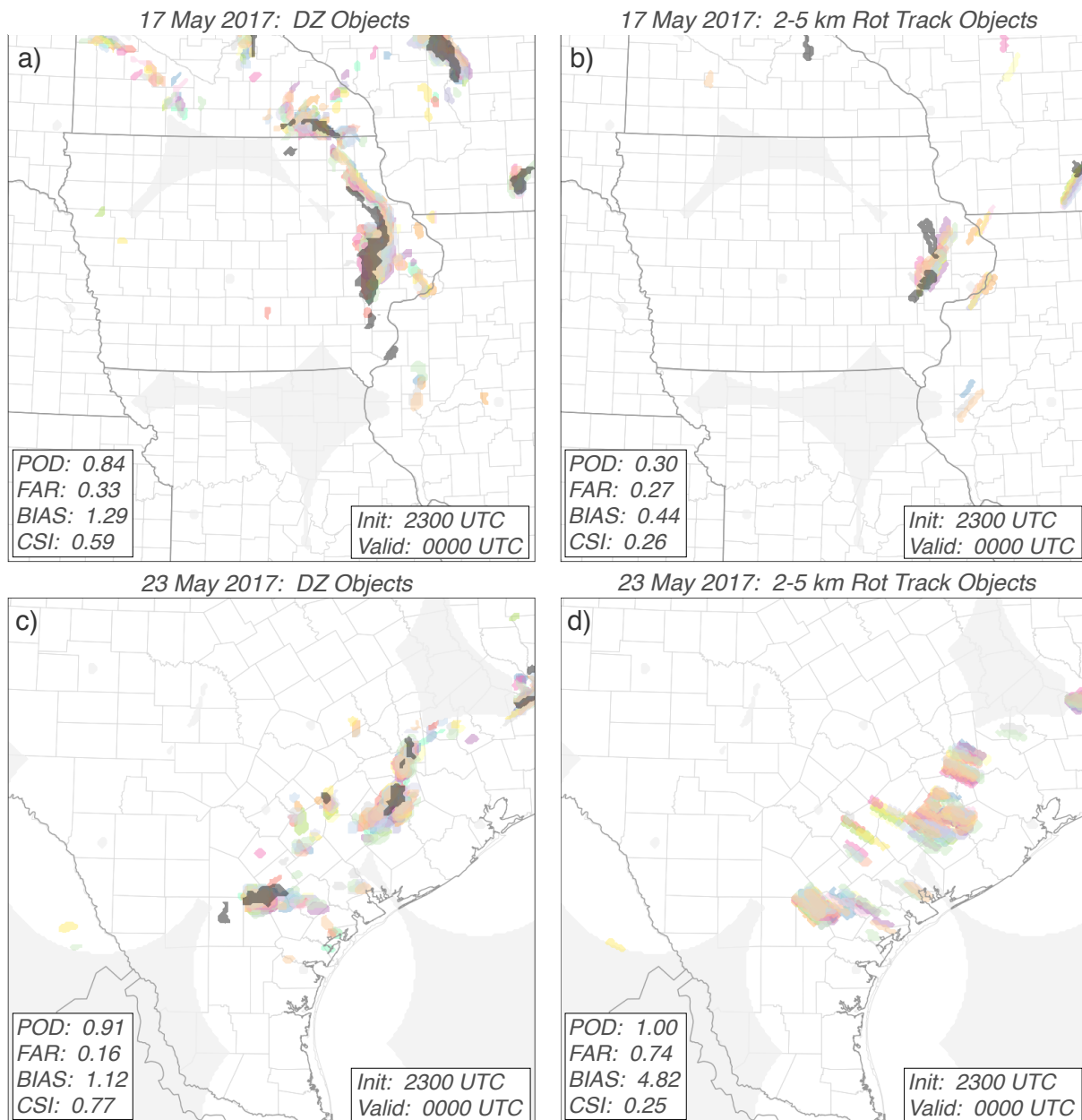


FIG. 14. As in Fig. 6 except for (a, c) composite reflectivity and (b, d) rotation track objects 60-minutes into forecasts initialized at 2300 UTC on (a, b) 17 May 2017 and (c, d) 23 May 2017. POD, FAR, BIAS, and CSI scores for each forecast are provided in the lower left of each panel. Note that some forecast rotation track objects in (d) are matched to observed objects at different times, resulting in a FAR less than 1.0 despite no observed objects being present at the forecast time plotted.

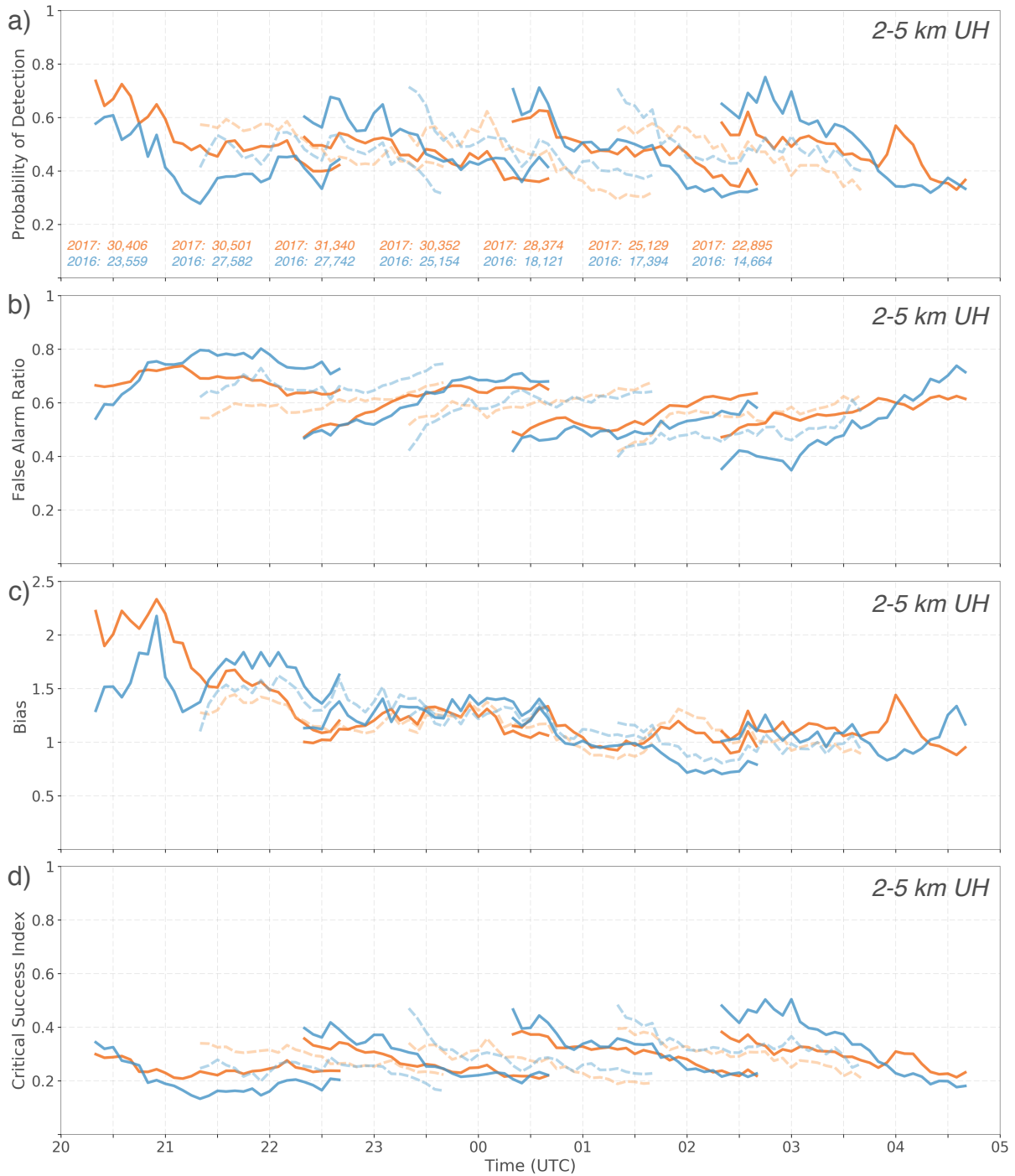


FIG. 15. As in Fig. 8 except for 2–5 km updraft helicity forecasts.



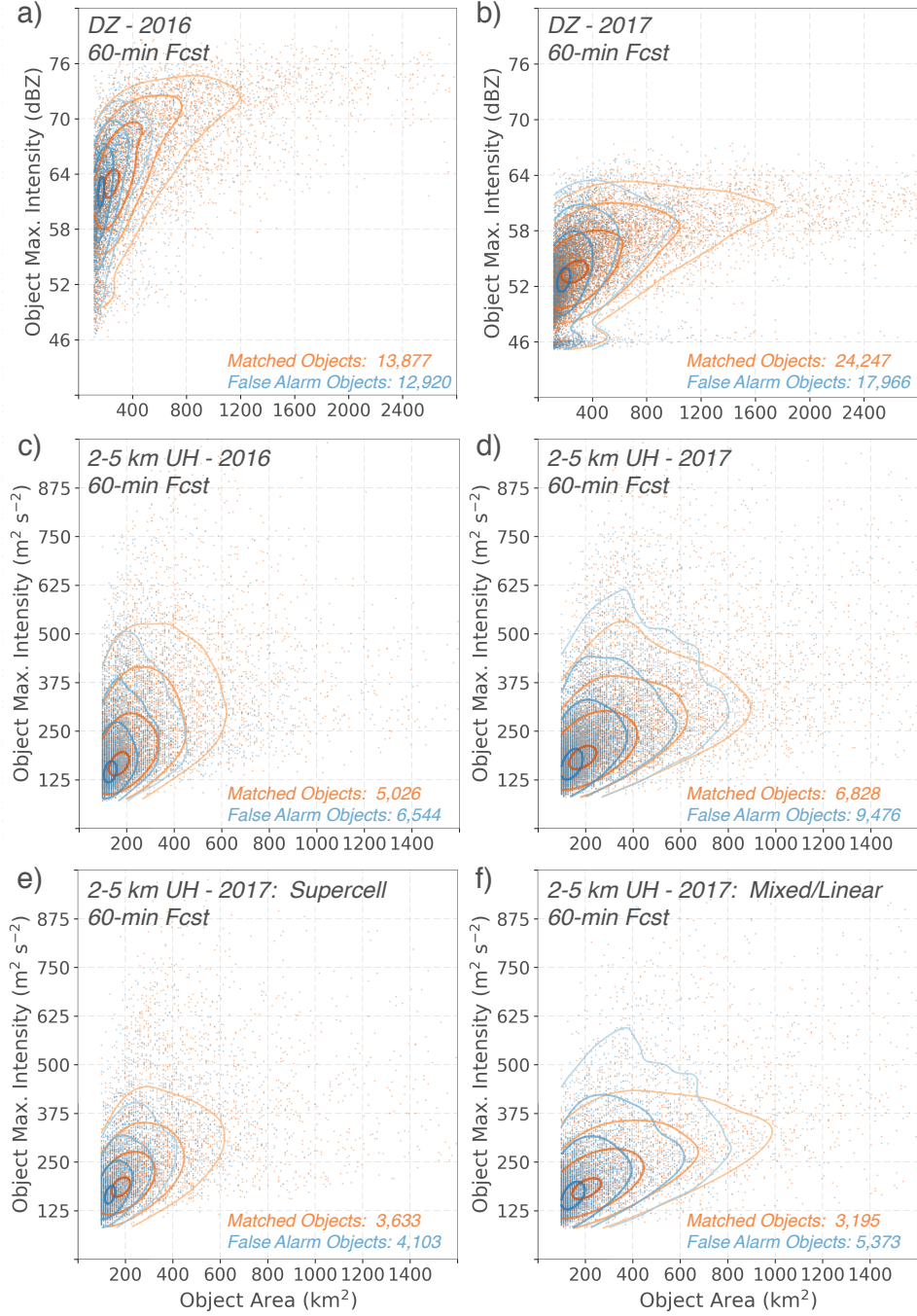


FIG. 16. Scatterplots of the parameter space of object area and maximum intensity for 60-minute NEWS-e forecasts of (a, b) composite reflectivity (dBZ) and (c–f) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ) during (a, c) 2016, (b, d) 2017, and 2017 cases classified as (e) supercell or (f) mixed/linear mode. Matched objects are plotted in orange and false alarm objects in blue with the total number of objects in each category listed in the lower right. Kernel density estimate contours of the 95th, 97.5th, 99th, and 99.9th percentile values of each distribution are overlain to illustrate differences between matched and false alarm distributions. Every third reflectivity object is plotted to improve clarity.

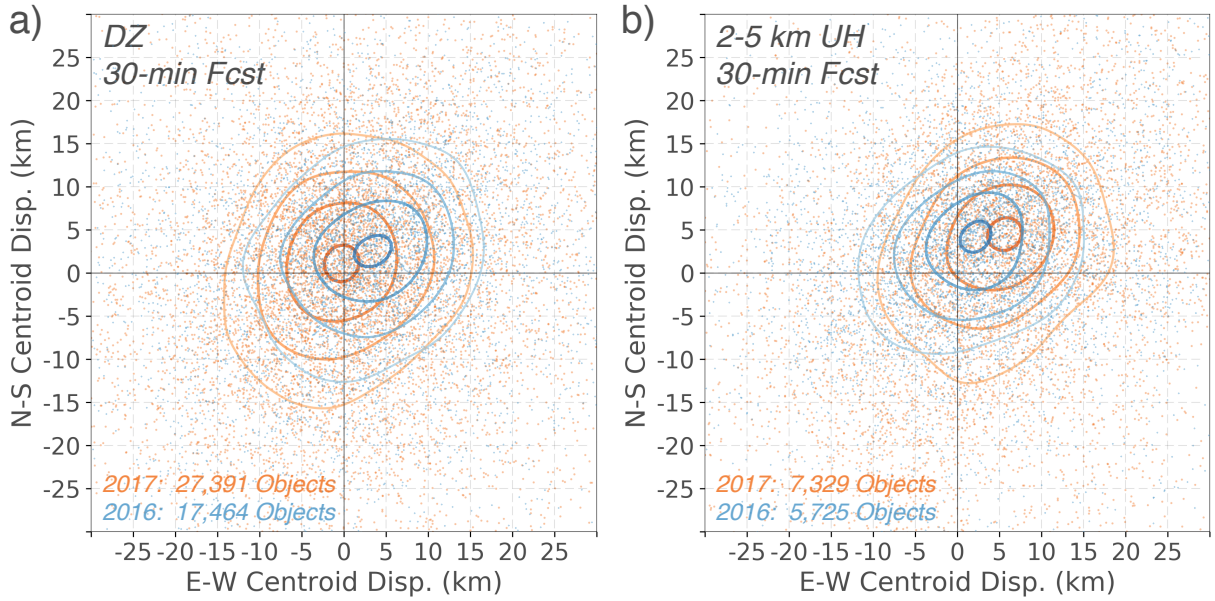


FIG. 17. Scatterplots of the east-west and north-south centroid displacements (km) of matched objects for 30-minute NEWS-e forecasts of (a) composite reflectivity (dBZ) and (b) 2–5 km updraft helicity ( $\text{m}^2 \text{s}^{-2}$ ). Objects from 2016 (2017) are plotted in blue (orange) and the total number of objects for each year is listed in the lower left. Kernel density estimate contours are overlain as in Fig. 16 and every third reflectivity object is plotted to improve clarity.

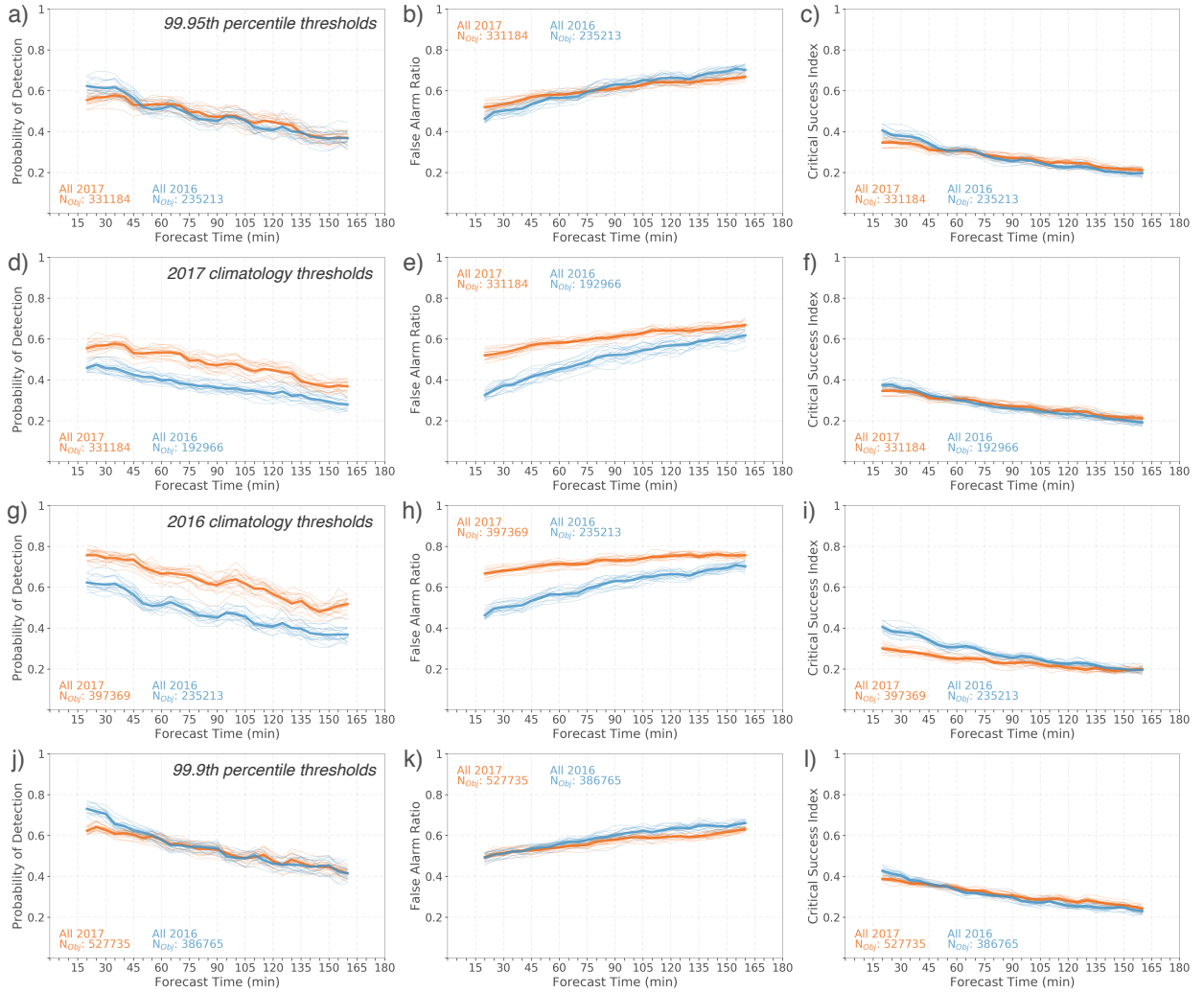


Fig. A1. Time series of (a, d, g, j) probability of detection, (b, e, h, k) false alarm ratio, and (c, f, i, l) critical success index for NEWS-e 2-5 km rotation track forecasts. The intensity threshold used to identify forecast and observed rotation track objects is varied between the (a–c) 99.95th percentile from each year’s climatology (same as Fig. 7), (d–f) the 99.95th percentile from the 2017 climatology only, (g–i) the 99.95th percentile from the 2016 climatology only, and (j–l) the 99.9th percentile from each year’s climatology. Individual ensemble member scores are plotted in thin orange (blue) lines with thick orange (blue) lines representing the ensemble mean for 2017 (2016) NEWS-e forecasts.

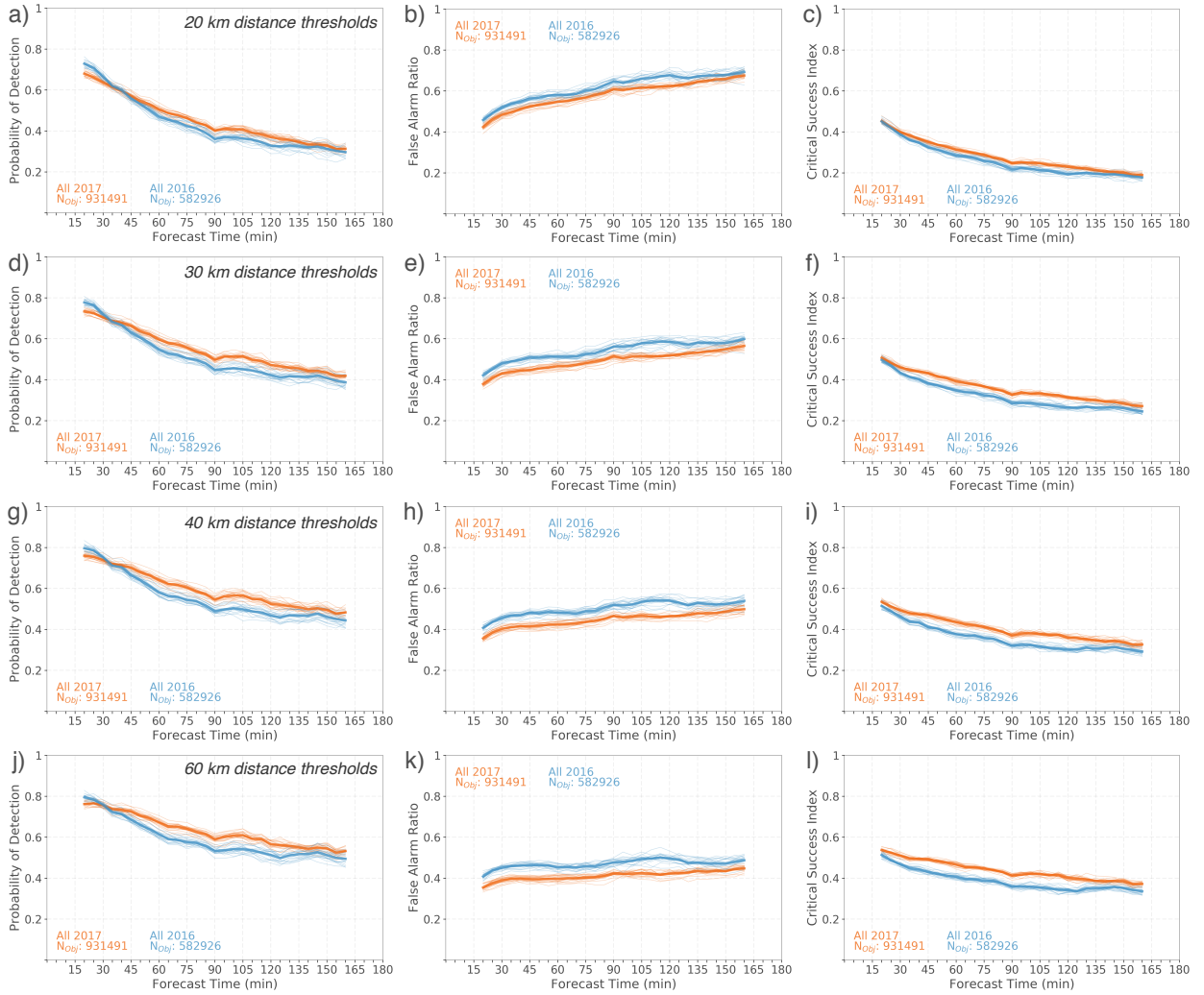


Fig. A2. As in Fig. A1, except for the composite reflectivity objects and the maximum distance threshold for object matching is varied from (a–c) 20 km, (d–f) 30 km, (g–i) 40 km (same as Fig. 5), and (j–l) 60 km.